THE OVERCOOKED GENERALISATION CHALLENGE

Anonymous authors

Paper under double-blind review

ABSTRACT

We introduce the Overcooked Generalisation Challenge (OGC) – the first benchmark to study reinforcement learning agents' zero-shot cooperation abilities when faced with novel partners and levels in the Overcooked-AI environment. This perspective starkly contrasts a large body of previous work that has evaluated cooperating agents only on the same level or with the same partner, thus failing to capture generalisation abilities essential for real-world human-AI cooperation. Our challenge interfaces with state-of-the-art dual curriculum design (DCD) methods to generate auto-curricula for training general agents in Overcooked. It is the first open-source cooperative multi-agent environment specially designed for DCD methods and, consequently, the first evaluated with state-of-the-art methods. It is fully GPU-accelerated, built on the DCD benchmark suite minimax, and freely available under an open-source license: http://anonymised.edu. We show that state-of-the-art DCD algorithms fail to produce useful policies on this novel challenge, even if combined with recent network architectures specifically designed for scalability and generalisability. As such, the OGC pushes the boundaries of real-world human-AI cooperation by enabling research on the impact of generalisation on cooperating agents.

025 026 027

028

000

001 002 003

004

006

008 009

010

011

012

013

014

015

016

017

018

019

021

023

1 INTRODUCTION

Developing computational agents capable of collaborating with humans has emerged as a key challenge in artificial intelligence (AI) research (Stone et al., 2010; Dafoe et al., 2020) and promises to vastly expand human abilities (O'neill et al., 2020). Recent years have seen considerable advances in understanding human cooperative behaviour (Rand & Nowak, 2013; Vizmathy et al., 2024), computational modelling of cooperation (Nikolaidis & Shah, 2013; Sadigh et al., 2016; Ding et al., 2024),as well as in developing computational methods for human-AI cooperation (Hu et al., 2020; Strouse et al., 2021). In parallel, several benchmarks (Samvelyan et al., 2019; Bard et al., 2020) were proposed to foster the development and evaluation of these methods. Most notably, Overcooked-AI (Carroll et al., 2019) has established itself as a widely used benchmark for evaluating (zero-shot) human-AI coordination (Strouse et al., 2021; Zhao et al., 2023; Yu et al., 2023).

Despite the advances they have enabled, all of these benchmarks are limited in that they only allow to assess reinforcement learning (RL) agents' 040 cooperative abilities in-distribution. That is, they either only allow to 041 evaluate agents in the same environment in which they were trained (Hu 042 et al., 2020; Carroll et al., 2019) or with the same partner agent they 043 were trained with (Foerster et al., 2018; Lowe et al., 2017; Strouse et al., 044 2021). In Overcooked-AI, for instance, existing zero-shot coordination (ZSC) methods are trained once per layout at considerable cost (Carroll et al., 2019; Yang et al., 2022; Zhao et al., 2023; Yu et al., 2023), 046 and these layouts only feature a limited number of possible cooperation 047 strategies (see Figure 1). However, real collaborative settings require co-048 ordination with novel partners in unknown environments. For example, consider a medical robot assisting doctors in hospitals. Such a robot will be deployed in unique and unknown hospitals and surgical rooms where 051 they need to adapt to different medical staff and their preferences. 052

053 To address this limitation, we introduce the *Overcooked Generalisation Challenge* (OGC) – the first zero-shot cooperation benchmark that chal-



Figure 1: Coordination challenges in the Overcooked-AI *Coordination Ring* layout.

lenges agents to cooperate in novel layouts and with unknown partner agents. While previous open-055 source benchmarks studied opponents combined with map generalisation, no dedicated open-source 056 benchmarks exist for studying cooperation partners combined with map generalisation. Cooperative 057 settings differ from competitive ones in their game-theoretic background and thus require separate 058 algorithms and benchmarks (Lerer & Peysakhovich, 2019). Neural MMO (Suarez et al., 2021; 2023) comes closest to our setting as it mixes cooperation and competition but crucially does not provide a purely cooperative setting which is specially designed for human-AI coordination – un-060 like Overcooked-AI. To train and evaluate agents on our benchmark, we make use of unsupervised 061 environment design (UED) (Dennis et al., 2020) to generate suitable training levels, provide hand-062 designed testing levels, and asses zero-shot cooperation on these by providing populations of diverse 063 testing agents. As such, our work is the first to combine UED techniques with a multi-agent RL zero-064 shot cooperation task and thus bridges the gap between two previously unrelated research areas; it 065 studies the impact of generalisation on human-AI coordination and the ability of UED algorithms to 066 design optimal auto-curricula for cooperating agents. We benchmark several UED algorithms and 067 network architectures on our challenge and find that they struggle to perform well. Only PAIRED 068 (Dennis et al., 2020), together with a policy that incorporates a soft Mixture-of-Experts (SoftMoE) 069 module (Obando-Ceron et al., 2024), has some limited success at generalising to the testing levels and outperforms competitive baselines, including robust PLR (Jiang et al., 2021b;a) and AC-070 CEL (Parker-Holder et al., 2022). Overall, our findings call for developing methods that combine 071 zero-shot coordination and DCD techniques in a single ZSC-DCD framework, and our benchmark 072 provides the environment to do so. Taken together, our contribution is three-fold: 073

- 074 1. We introduce the Overcooked Generalisation Challenge – a novel benchmark challenge in which 075 agents are asked to cooperate with novel partners in previously unseen layouts.
 - 2. We provide OvercookedUED an open-source environment that can be used with state-of-the-art DCD algorithms and that is integrated into minimax (Jiang et al., 2023), taking full advantage of the hardware acceleration provided by JAX.
 - 3. We benchmark our environment by training agents with common DCD algorithms (Dennis et al., 2020; Jiang et al., 2021a; Parker-Holder et al., 2022) and show that current DCD algorithms struggle with the challenge even if we employ recent network architectures (Smith et al., 2023; Obando-Ceron et al., 2024). Furthermore, we assess zero-shot cooperation performance with a population of diverse partners to link zero-shot cooperation and generalisation. We show that as policies become more generally capable, they achieve better zero-shot cooperation.
- 085 087

088

090 091

076

077

078

079 080

081

082

083

084

2 **RESEARCH CHALLENGES**

The OGC poses several new challenges for zero-shot human-AI cooperation that go beyond existing benchmarks and that are essential for further advances in the development of cooperating RL agents:

Generalisation The OGC challenges the generalisation capabilities of methods and agents by 092 having them engage in a double generalisation challenge: adjusting to both novel partners and lev-093 els. Existing cooperative open-source benchmarks require typically only one form of generalisation, see for instance (Lowe et al., 2017; Foerster et al., 2018; Carroll et al., 2019; Hu et al., 2020). 095

096

094

Environment Design Our environment challenges UED algorithms in generating and designing 097 layouts with many interacting components and agents. This is in contrast to existing environments 098 that only require UED algorithms to design simple mazes, 2D walker terrains, or race tracks consisting of fewer elements (Dennis et al., 2020; Jiang et al., 2021a; Parker-Holder et al., 2022; 100 Rutherford et al., 2024a). We show that methods struggle to design layouts similar to the ones 101 humans designed. Current methods specifically fail to design layouts requiring handing over 102 items over a countertop or featuring deliberately designed circuits. Our benchmark challenges 103 further research to develop UED methods that design more realistic collaboration environments for 104 curriculum learning, possibly along the lines of Bruce et al. (2024).

105

Combining Environment and Partner Generalisation Coordinating with novel partners and 106 generalising to novel levels were often treated as separate research areas. As such population-based 107 methods for zero-shot coordination do not apply to the challenge since training levels are generated

111 112 Name Multi-Zero-GPU Open Partial Img. obs. agent shot accel-Source obs. 113 erated coop. 114 115 ? XLand (Team et al., 2021; Bauer et al., 2023) \checkmark 116 LaserTag (Samvelyan et al., 2023) MultiCarRacing (Samvelyan et al., 2023) 117 \checkmark 118 CoinRun (Cobbe et al., 2019) ~ √ -119 ProcGen (Cobbe et al., 2020) -**√** √ *√* -120 2D Mazes (Cobbe et al., 2019; Dennis et al., 2020) √ 121 CarRacing (Jiang et al., 2021a) _ √

-

√

 \checkmark

 \checkmark

 \checkmark

 \checkmark

./

 \checkmark

√

√

 \checkmark

 \checkmark

 \checkmark

Table 1: Overview of benchmarks for unsupervised environment design and procedurally generated
 environments. Closed-source benchmarks are marked in gray – these cannot be evaluated on by the
 research community.

 \checkmark

 \checkmark

1

-

 \checkmark

134 135

136 137

138

122

123

124

125

126 127

128 129

3 RELATED WORK

3.1 GENERALISATION IN REINFORCEMENT LEARNING

Bipedal Walker (Wang et al., 2019)

XLand-MiniGrid (Nikulin et al., 2023)

AMaze (Jiang et al., 2023)

Craftax (Matthews et al., 2024)

JaxNav (Rutherford et al., 2024a)

OvercookedUED (ours)

139 A large number of works have shown that RL agents fail to generalise to new environments, 140 see (Zhang et al., 2018a; Cobbe et al., 2019), and have triggered research on the generalisation 141 capabilities of RL agents (Nichol et al., 2018; Cobbe et al., 2019; 2020). Early results revealed that 142 RL agents can memorise large numbers of levels during training (Zhang et al., 2018b; Cobbe et al., 2019) and that they must experience sufficiently diverse training data to generalise well (Cobbe 143 et al., 2020). One established approach to generate diverse training data is to use domain ran-144 domisation (Jakobi, 1997, DR). Still, DR has been shown to produce many uninformative samples 145 (Khirodkar et al., 2018), which can lead to the agent's inability to generalise (Dennis et al., 2020). 146

147 148

3.2 UNSUPERVISED ENVIRONMENT DESIGN

149 Intending to address this challenge, later works on generalisation focused on unsupervised environ-150 ment design (Dennis et al., 2020, UED). UED aims to improve domain randomisation by generat-151 ing auto-curricula that include training levels of increasing complexity to facilitate continued agent 152 learning (Graves et al., 2017). It does so by adapting the free parameters of an under-specified en-153 vironment to the agent's capabilities. Most popular UED methods fall into the category of Dual 154 Curriculum Design (Jiang et al., 2021a, DCD) that combine 1) an agent, 2) a level generator, and 155 3) a curator that picks which levels to train on. Popular methods include Prioritised Level Replay (PLR) (Jiang et al., 2021b), robust PLR[⊥] (Jiang et al., 2021a), MAESTRO (Samvelyan et al., 156 2023), ReMiDi (Beukman et al., 2024), PAIRED (Dennis et al., 2020), ACCEL (Parker-Holder et al., 157 2022), and Replay-Enhanced (RE)PAIRED (Jiang et al., 2021a). While the development of these 158 DCD methods has been steady, they have mostly been explored in simple environments, see Table 1. 159

- 160
- 161 **Single-agent UED Environments** Early work on generalisation mainly focused on single-agent environments (Zhang et al., 2018b; Farebrother et al., 2018; Cobbe et al., 2019) and these are also

<sup>on the fly. Thus, training a best response against a diverse population on each layout is infeasible.
There is currently no algorithm to train a population of diverse agents over a distribution of levels.
Our benchmark encourages these branches to merge both lines of research and develop UED-ZSC methods, i.e. methods that learn both at the same time.</sup>

popular in UED research. Among these, prior work has studied mazes (Dennis et al., 2020; Jiang et al., 2021a; Parker-Holder et al., 2022; Jiang et al., 2023; Li et al., 2023a; Beukman et al., 2024), bipedal walkers (Wang et al., 2019; 2020; Parker-Holder et al., 2022) or car racing environments (Jiang et al., 2021a). One likely reason for their popularity as benchmarks for DCD is that new levels are easy to generate, and agents are usually fast to train. However, they are limited to a single agent, with limited options to interact with the environment and other agents, and thus bear little resemblance to real-world problems.

169

170 Multi-agent UED Environments Compared to single-agent environments, multi-agent environments are inherently more complex because the agents interact with each other, as well as 171 with the physical environment. Multi-agent environments are still rarely used in UED research. 172 Most prominent is Deepmind's XLand (Team et al., 2021; Bauer et al., 2023), a closed-source 173 multi-task universe for generating single- and multi-agent tasks and environments. While XLand 174 features cooperative tasks, it is not available to researchers for studying cooperative multi-agent 175 UED algorithms. While an open-source variant was recently published (Nikulin et al., 2023), it 176 only supports a single agent. Arguably closest is Neural MMO (Suarez et al., 2021; 2023), which 177 is a massively multi-task and multi-agent environment that mixes cooperation and competition to 178 replicate massively multiplayer online games. We instead are interested in assessing and identifying 179 cooperation performance in specially designed human-AI cooperation challenges for which the 180 maissvely multi-task and multi-agent cooperation-competition setting of Neural MMO is unsuitable. 181 Additionally, classic Overcooked already benefits from a rich history of human-AI cooperation research. Finally, while JaxNav (Rutherford et al., 2024a) features multi-agent path-finding no 182 interaction between agents is required and the environment is not focused on human-AI cooperation. 183 Other open-source environments are competitive, i.e. LaserTag (Lanctot et al., 2017; Samvelyan 184 et al., 2023) and MultiCarRacing (Schwarting et al., 2021; Samvelyan et al., 2023), and thus not 185 applicable to our setting. Opposed to all of these, our work contributes to and analyses the first open-source cooperative multi-agent UED environment.

187 188

3.3 HUMAN-AI COOPERATION IN OVERCOOKED-AI

190 Overcooked-AI (Carroll et al., 2019) has become one of the most important benchmarks for 191 human-AI cooperation. The environment is fully cooperative and has two agents cook and deliver 192 soups to earn a joint reward. Overcooked-AI was, for example, used in research on zero-shot 193 cooperation (Strouse et al., 2021; Zhao et al., 2023; Yu et al., 2023; Li et al., 2023b; Yan et al., 194 2023, ZSC), language model-based cooperative agents (Liu et al., 2024; Tan et al., 2024), human 195 modelling in cooperation (Yang et al., 2022). Zero-shot cooperation refers to cooperating with a 196 partner not encountered during training. It is an important proxy to ensure the ability of an agent to coordinate with humans at test time, given that human data is often unavailable and agents thus must 197 be able to coordinate effectively without previous training. It is commonly studied in Overcooked. 198

Related to our work is the work of Fontaine et al. (2021) in which the authors used procedurally generated Overcooked layouts to evaluate the impact of different layouts on human-robot interaction using planning algorithms. However, while they use procedural content generation in the Overcooked context their research does not focus on cross-layout generalisation – a major theme in our work. Our work is thus the first to explore the impact of cross-level generalisation for zero-shot cooperation and is the first to provide the necessary tools for this.

205 206

207

4 PRELIMINARIES

208 The cooperative multi-agent UED setting can be formalised as a decentralised under-specified par-209 *tially observable Markov decision process* (Dec-UPOMDP) with shared rewards. A Dec-UPOMDP is defined as $\mathcal{M} = \langle \mathcal{N}, A, \Omega, \Theta, \mathcal{S}^{\mathcal{M}}, \mathcal{T}^{\mathcal{M}}, O^{\mathcal{M}}, \mathcal{R}^{\mathcal{M}}, \gamma \rangle$ in which \mathcal{N} is the set of agents with 210 211 cardinality n, Ω is a set of observations, and $\mathcal{S}^{\mathcal{M}}$ is the set of true states in the environment. Partial 212 observations $o^i \in \Omega$ are obtained by agent $i \in \mathcal{N}$ using the observation function $O: \mathcal{S} \times \mathcal{N} \to \Omega$. 213 Following Jiang et al. (2021a), a *level* \mathcal{M}_{θ} is defined as a fully-specified environment given some parameters $\theta \in \Theta$. In it, agents each pick an action $a_i \in A$ simultaneously to produce a joint action 214 $\boldsymbol{a} = (a_1, \ldots, a_n)$ and observe a shared immediate reward $R(s, \boldsymbol{a})$. Then, the environment transi-215 tions to the next state according to a transition function $\mathcal{T}: \mathcal{S} \times \mathcal{A}^1 \times ... \times \mathcal{A}^n \times \Theta \to \Delta(\mathcal{S})$ where



Figure 2: Overview of the Overcooked Generalisation Challenge (OGC) and how it is typically used in a Dual Curriculum Design (DCD) algorithm. The OGC supports teacher-based methods like PAIRED (Dennis et al., 2020) via unsupervised environment design (UED) and edit-based methods like ACCEL (Parker-Holder et al., 2022) via mutator functions of existing layouts.

234 $\Delta(S)$ refers to the space of distributions over S. $\gamma \in [0,1)$ specifies the discount factor. Agents learn a policy π . The joint policy π together with the discounted return $R_t = \sum_{i=0}^{\infty} \gamma^i r_{t+1}$ induce 235 a joint action value function $Q^{\pi} = \mathbb{E}_{s_{t+1:\infty}, \boldsymbol{a}_{t+1:\infty}}[R_t|s_t, \boldsymbol{a}_t]$. The definition of the Dec-UPOMDP extends a Dec-POMDP (Oliehoek & Amato, 2016; Wu et al., 2021) with the free parameters of the 236 237 environment Θ , analogously to previous works (Dennis et al., 2020; Jiang et al., 2021a; Samvelyan 238 et al., 2023). Our definition differs from (Samvelyan et al., 2023) in terms of the shared rewards 239 and general-sum nature. Within our Dec-UPOMDP, we perform UED to train a policy over a distri-240 bution of fully specified environments that enable optimal learning. This is facilitated by obtaining 241 an *environment policy* Λ (Dennis et al., 2020) that specifies a sequence of environment parameters 242 Θ^{I} for the given policy that is to be trained. How Λ is obtained depends on the DCD method. 243 For example, in OvercookedUED, Θ represents the possible positions of walls, pots, serving spots, 244 agent starting locations, and onion and bowl piles which is adjusted by Λ throughout training.

245 246

228

229

230

231 232 233

5 THE OVERCOOKED GENERALISATION CHALLENGE

247 248

An overview of the Overcooked Generalisation Challenge is shown in Figure 2. The OGC extends 249 previous work by evaluating the cooperative abilities out-of-distribution. That is, in contrast to 250 existing UED environments, the OGC focuses on the cooperation of multiple agents in a complex, 251 cooperative task across different levels. More specifically, two different agents are tasked with 252 cooking a soup together in the five original layouts of Overcooked-AI (see Figure 3), but without 253 having encountered them and their partner during training. Since the original five layouts have been 254 designed to test and explore different kinds of cooperation, they form suitable out-of-distribution 255 test levels. To train an agent capable of generalisation, we generate a curriculum of possibly endless 256 diverse training layouts via procedural content generation. The OGC is more closely aligned 257 with real-world human-AI collaboration as it does not limit evaluation to one specific physical 258 environment or partner. To generate a curriculum of layouts, we use DCD methods. Specifically, 259 methods in which an environment designer interacts with the challenge by designing layouts 260 either from scratch through interacting with *OvercookedUED* – a novel environment for creating Overcooked levels – by alternating existing layouts through the Overcooked mutator or by letting 261 the OGC generate random layouts. At every step of the curriculum, this designer must account for 262 agents' cooperation ability when trying to generate layouts that are at the forefront of their abilities. 263

264 265

266

5.1 COMPONENTS OF THE CHALLENGE

While OGC refers to the challenge as a whole, it comprises several components that enable its integration with DCD algorithms (see Figure 2). At the heart of it, it features an Overcooked environment capable of running different levels fast and in parallel in which agents learn to cooperate. It features *OvercookedUED* that provide methods, interfaces and a teacher environment

278

279 280

281

282 283

295



Figure 3: We study the five layouts proposed by Carroll et al. (2019). From left to right: *Cramped Room, Asymmetric Advantages, Coordination Ring, Forced Coordination, and Counter Circuit.*

to design novel layouts as well as an *Overcooked Mutator* that alters existing layouts, specifically designed to be used with ACCEL.

Overcooked-AI OGC builds on the Overcooked-AI environment. We adapted the version 284 provided by the JaxMARL project (Rutherford et al., 2024b), keeping features consistent with the 285 original implementation. This includes action and observation spaces, i.e. the set of actions is 286 {left, right, up, down, interact, stay} and observations are encoded as a stack of 26 287 $\dot{h} \times w$ boolean masks encoding the positions of elements in the environment. In this representation, 288 the first mask encodes the position of the first agent, the second mask the one of the second agent 289 etc. Since agents now learn to play on many different layouts all at once, we adjust the environment 290 to be capable of parallelising across differently shaped levels via padding. I.e., during rollouts, 291 layouts are padded to a maximum size, and all objects in these layouts are one-hot encoded based 292 on their position in equally sized masks. While this facilitates fast parallel rollouts that can be 293 just-in-time compiled, it requires the introduction of a maximum height h and width w that need to be picked as a hyperparameter before training. 294

296 **OvercookedUED** OvercookedUED features the interfaces necessary to design new layouts. For algorithms that make use of a teacher agent to create layouts (PAIRED, etc.), OvercookedUED 297 provides a teacher environment (see Figure 2). This environment allows a teacher policy to take 298 design steps to parameterise the underspecified MDP. At every timestep t of the generation process 299 the teacher observes the unfinished layout and picks an action from a space that consists of the total 300 number of cells in the $h \times w$ grid. This cell then becomes filled with the next items to be placed. 301 Objects are placed sequentially and in a deterministic order, starting from walls, agents one and two, 302 goal, onion, pot and bowl positions. An episode in the teacher MDP lasts until all items are placed. 303 In case of a conflict, elements are placed randomly on free cells. The teacher is parameterised by 304 its own neural network. As in previous work (Jiang et al., 2023), OvercookedUED does not check 305 whether a layout is solvable and leaves the task of designing and/or identifying suitable training 306 layouts to the DCD method.

For algorithms that do not specify a teacher, such as PLR, OvercookedUED generates random layouts. These random layouts feature one or two piles of onions, bowls, pots and serving locations, and two agents.

Finally, some DCD algorithms, such as ACCEL, require alternating existing layouts by mutating them. OvercookedUED supports layout mutation through five basic operations: (1) converting a random wall to a free space and vice versa, (2) moving goals, (3) pots, (4) plate piles, and (5) onion piles. Given a layout, our *mutator* randomly samples *n* operations and applies them. All versions allow the number of walls placed to be configured and the environment always places a border wall.

316 Implementation The OGC is implemented in Jax (Bradbury et al., 2018) and integrated into 317 minimax (Jiang et al., 2023). As such, it can be tested with all available DCD algorithms present 318 in minimax. To achieve this we extend minimax with runners, replay buffers etc. that are 319 compatible with multiple agents. Building on an established library eliminates sources of error and 320 presents users of the challenge with a familiar experience. We present the steps-per-seconds (SPS) 321 on our setup given varying degrees of parallelism in Table 2 and compare it to the GPU-accelerated maze environment minimax includes AMaze. Given sufficiently large numbers of parallel 322 environments, OGC can be run at hundreds of thousands of SPS. While less than AMaze, the OGC 323 is a more fully-featured environment in which multiple agents take steps and interact.



Figure 4: Sample levels generated by the different methods after 15,000 (Middle) and 30,000 (End) epochs. Even after considerable training, none of the methods can guarantee the generation of solvable layouts (Middle-row leftmost and rightmost).

5.2 EVALUATION

343 We evaluate agents by their per-344 on out-of-distribution formance Overcooked-AI layouts to asses gen-345 eralisation performance in self-play 346 and in cross-play. In cross-play, 347 a fictitious co-play (Strouse et al., 348 2021, FCP) and maximum entropy-349 based population based training 350 (Zhao et al., 2023, MEP) population 351 of a total 24 agents each is used to 352 asses zero-shot cooperation. Both

Table 2: Average steps-per-second for different numbers of parallel environments measured by taking 1,000 steps with randomly sampled actions.

# Parallel Envs	1	32	256	1024
AMaze OvercookedUED	$264 \\ 151$	${8,141 \atop 4,921}$	$\begin{array}{c} 67,282 \\ 40,011 \end{array}$	$\begin{array}{c} 264,142 \\ 156,696 \end{array}$

353 populations are trained with equal settings and include a low, medium and high-skilled checkpoint 354 of each run extracted at 10, 50 and 100 % achieved return respectively. The population entropy 355 coefficient α is 0.01 for MEP. In this work we define zero-shot coordination as the task of cooperating with a partner previously not encountered during training and view it in contrast to 356 ad-hoc teamwork (Stone et al., 2010) since in our setting there is no time to update a fixed policy 357 after training (Hu et al., 2020). As zero-shot cooperation with a diverse population has become a 358 proxy for assessing the abilities of an agent to coordinate with humans. Our benchmark includes 359 the necessary tools to perform this evaluation. In our analysis, we report results using the mean 360 episode reward and mean layout solved rate, similar to previous work (Jiang et al., 2023). A layout 361 is considered solved if an agent pair delivers more than one soup which differentiates goal-directed 362 from random behaviour. We present these metrics in the self- and cross-play settings. Additionally, 363 we investigate what kinds of levels agents perform poorly in and why in a final error-analysis.

364 365 366

367

6 ANALYSING & BENCHMARKING THE CHALLENGE

We benchmark the challenge with several DCD algorithms and network architectures. We aim to set a performance baseline for future works and show what evaluations this benchmark enables. To this end, we first show that generalising to novel layouts in Overcooked is difficult, and then we move on to the additional challenge of zero-shot cooperation.

All baselines are trained using MAPPO, which is known to work well in cooperative settings (Yu et al., 2022) using centralised training and decentralised execution (Foerster et al., 2016). As for DCD algorithms, we compare the performance of DR, PLR^{⊥, ||}, Pop. PAIRED and ACCEL^{||}. We chose these methods as they have better theoretical guarantees (PLR[⊥] vs PLR), better runtime performance (ACCEL^{||} and PLR^{||}), or because we found them to perform better empirically (Pop. PAIRED vs PAIRED). We excluded POET (Wang et al., 2019) in this analysis as it outputs specialists rather than generalists, which we require (Parker-Holder et al., 2022). Additionally, we excluded

Table 3: Mean episode reward for the different methods averaged over the respective testing layouts.
The best result is shown in **bold**. We report aggregate statistics over three random seeds. We include
Oracles which were trained on the five testing layouts directly to establish a empirical maximum.

Method	CNN-LSTM	SoftMoE-LSTM	CNN-S5
DR	0.46 ± 0.16	5.22 ± 7.19	0.00 ± 0.00
$PLR^{\perp,\parallel}$	0.17 ± 0.06	0.91 ± 0.71	0.12 ± 0.13
Pop. PAIRED	0.19 ± 0.09	13.34 ± 5.70	0.24 ± 0.19
ACCEL∥	0.20 ± 0.14	0.67 ± 0.60	0.28 ± 0.2
Oracle	189.49 ± 12.96	217.02 ± 39.18	$\textbf{155.01} \pm 12.8$

388 389 390

381 382

391 MAESTRO as it is based on prioritised fictitious self-play (Heinrich et al., 2015; Vinyals et al., 2019) 392 that is not easily adaptable to the cooperative setting (Strouse et al., 2021). As in (Jiang et al., 2023), if not stated otherwise, we train in 32 parallel environments and stop after 30,000 outer training 393 loops, amounting to just under 400 million steps in the environment. Hyperparameters were picked 394 after a grid search over reasonable values, and all parameters are provided in Appendix A.4. Our de-395 fault neural network architecture consists of a convolutional encoder with a recurrent neural network 396 with an LSTM (Hochreiter & Schmidhuber, 1997). It is picked for its good performance in previous 397 work (Yu et al., 2023) (see Appendix A.5 for details). In addition to our default network architecture, 398 we explore the use of SoftMoE (Obando-Ceron et al., 2024), which have recently been identified 399 for their potential for enabling scaling and generalisation, and S5 layers (Smith et al., 2023) due to 400 the strong results of structured state-space models (Gu et al., 2022) in meta reinforcement learning 401 (Lu et al., 2023). SoftMoE modules replace the penultimate layer after the feature extractor and S5 402 layers the LSTM in all experiments. We hypothesise that these provide better generalisation performance. Using these parameters, we verified that agents also overfit to their level in Overcooked by 403 evaluating agents trained on a single layout on all layouts (cf. Appendix A.6.1). Additionally, we 404 verified that all architectures can be fitted to the testing layouts when trained on them directly. We 405 will refer to these as *Oracles* and use them to establish the maximum performance possible. Lastly, 406 for all runs we display training curves on seen and the five unseen evaluation levels in Appendix 407 A.7.1. 408

409

Layout Generalisation Performance Simply generalising to the testing layouts in the OGC is 410 already challenging for all methods without having to coordinate with novel partners, as presented 411 in Table 3. Compared to commonly used single-agent Maze environments (such as AMaze, 412 compare (Jiang et al., 2023)), all DCD methods struggle to obtain good results. This is most 413 evident when compared with oracle policies (bottom row). PAIRED outperforms all other models 414 significantly 0.01 using a one-sided paired t-test. This is also shown in the mean415 solved rate where it reaches $14.6 \pm 7.7\%$, while all other models have a solved rate of mostly 0%416 (cf. Appendix A.6.2). While this model performs better on average, layouts differ greatly in their 417 difficulty. Our best-performing model reaches modest performance in Asymmetric Advantages 418 and Cramped Room while mostly failing in the others, with no other model achieving noteworthy 419 results. Recall that the environment features more moving parts that must be placed correctly to 420 facilitate learning. This makes it hard for approaches like DR to find optimal placements by pure 421 chance, as reflected in the results. The full results are in the Appendix A.6.3.

422

423 **Zero-Shot Cooperation Performance** Ultimately, we want the OGC to connect map generalisa-424 tion and zero-shot coordination. To that end, we train and then use a FCP and MEP population (see 425 Appendix A.6.4 for details) to establish how general cooperative agents can coordinate with diverse 426 policies. We present preliminary results in Figure 5 together with two other baselines: stay which is 427 a partner that never moves and *random* which samples random actions. As performance on out-of-428 distribution levels rises, agents become more competent at zero-shot cooperation. PAIRED always 429 outperforms baselines (cf. Appendix A.6.5). However, even PAIRED policies often perform only slightly better than random baselines, which signifies the challenges of our benchmark. This is also 430 evidenced by the kinds of levels these methods generate (Figure 4), as they tend to pivot towards gen-431 erating open spaces that ease cooperation but are notably different from evaluation layouts. Overall,

464

465

466

467 468 469

470 471



Figure 5: Zero-shot coordination results of the SoftMoE-LSTM policy paired with an FCP population trained on the respective layout. We report the mean episode reward and standard error.



Figure 7: Sample levels that our models perform best (top) and worst (bottom) in. The number of visits to each grid cell is shown as a heatmap overlay, while the mean return is stated below each layout. The layouts the model performs worst in tend to feature narrow elements or large distances between items.

cooperation performance is mostly carried by the expert FCP and MEP agents (compare Tables 13 and 14), mostly since our agents struggle to perform on the evaluation layouts in the first place.

Error Analysis We perform two final experiments to investi-472 gate the poor performance of our baselines and to eliminate trivial 473 sources of errors. First, we hypothesise that the top-down observa-474 tions in OGC are hard to generalise from since they are not invariant 475 to mirroring or rotations (Ye et al., 2020). To test this we evaluate 476 agents on 24 hand-designed circular evaluation levels with different 477 kinds of symmetry, as shown in Figure 6. We find that agents tend 478 to perform similarly across these layouts as the standard deviation 479 is at most 1.1, and therefore reject this hypothesis (more details in 480 Appendix A.7). Second, we investigate the kinds of levels our best-481 performing model does well vs poorly in from a pool of randomly 482 generated evaluation levels in Figure 7. The figure summarises the 483 cooperation behaviour of the agents by showing which cells are visited most frequently to give an impression of their motion patterns. 484 While on many layouts our model reaches good self-play perfor-485 mance (up to a maximum mean reward of 84.4; top row), it typically



Figure 6: An illustration of the circular evaluation levels; we move the kitchen around the sides and vary the size.

delivers few to no soups in layouts it performs worst in. These levels tend to be narrow/convoluted
and/or feature big distances between objects. Notice that the training levels in which our model
performs well in are similar to Asymmetric Advantages and Cramped Room, while the worst levels are similar to the other 3 evaluation levels. In conclusion, current DCD methods struggle with
generating training layouts of the correct complexity, i.e. ones that are similarly hard to evaluation ones.

Discussion Previous work (Jiang et al., 2021a) has found that PLR^{\perp} tends to outperform the 493 other here-tested algorithms in navigation-based tasks. Our more challenging environment suggests 494 that this might not always be the case. In our preliminary analysis, PAIRED outperformed other 495 DCD methods. Compared to mazes, car racing, or walker environments with fewer moving pieces, 496 Overcooked layouts are more complex to design, requiring the designer to place multiple objects 497 in relation to each other and the agents. Methods that employ a random generator therefore 498 struggle in such a big design space. This thus requires a capable generator and suggests that 499 simple navigation-based environments used to benchmark DCD in UED algorithms do not allow 500 full performance evaluation. As such, OvercookedUED can be an important part of evaluating 501 DCD algorithms. We envision that general Overcooked agents should be evaluated in scenarios 502 that are difficult for self-play agents using our benchmark. These include zero-shot cooperation 503 with strongly-biased agents (Yu et al., 2023) in Coordination Ring (see Section 1) and Asymmetric Advantages as described in (Ruhdorfer, 2023) and for which we provide the tools. 504

505 506

507

492

7 LIMITATIONS

508 Despite its many advantages, our challenge has two limitations. First, we artificially restricted the 509 maximum size of the layouts to allow the environment to be both fully observable as in Carroll et al. (2019) and parseable by CNN-based feature encoders. Future work should focus on more 510 natural representations of the whole scene, e.g. using graphs or item embeddings. While we 511 included a partial observation that could theoretically be computed independently of size, similar 512 to the vector-based observation used for behaviour cloning agents in (Carroll et al., 2019), batching 513 across layouts in OvercookedUED still requires the layouts to be scaled to the same height and 514 width. Second, while our challenge allows us to study zero-shot coordination via generalising 515 across layouts, reasoning about other agents (Rabinowitz et al., 2018; Gandhi et al., 2021; Bara 516 et al., 2023; Bortoletto et al., 2024b;a) might be equally important to achieve zero-shot cooperation 517 capabilities on unknown layouts. This is plausible given that humans can reason about the mental 518 states of other agents via Theory of Mind (Premack & Woodruff, 1978), as well as the physical 519 configuration of the space in which they operate. Future work could thus explore reasoning about 520 other agents in previously unexplored environments.

521 522

523

533

8 CONCLUSION

524 We have presented the Overcooked Generalisation Challenge (OGC) – a generalisation challenge 525 focusing on (zero-shot) cooperation in MARL in out-of-distribution test levels. Our challenge is the 526 first open-source cooperative multi-agent UED environment and is significantly more challenging 527 than previous environments commonly used in UED and DCD research. In addition to using 528 the challenge in UED research, we have shown how the OGC can be used in future research on human-AI collaboration as a zero-shot cooperation benchmark for general agents. That is, our 529 challenge establishes a link between generalisation and zero-shot coordination. Our work is the 530 first to provide the research community with the tools to train and evaluate agents capable of 531 coordinating in previously unknown physical spaces and with novel partners. 532

- 534 REFERENCES
- Lei Jimmy Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization. CoRR, abs/1607.06450, 2016. URL http://arxiv.org/abs/1607.06450.
- Cristian-Paul Bara, Ziqiao Ma, Yingzhuo Yu, Julie Shah, and Joyce Chai. Towards collaborative
 plan acquisition through theory of mind modeling in situated dialogue. In *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence*, pp. 2958–2966, 2023.

564

565

566

567

568 569

570

571

572

573

580

581

582

583

584

585

591

Nolan Bard, Jakob N. Foerster, Sarath Chandar, Neil Burch, Marc Lanctot, H. Francis Song, Emilio Parisotto, Vincent Dumoulin, Subhodeep Moitra, Edward Hughes, Iain Dunning, Shibl Mourad, Hugo Larochelle, Marc G. Bellemare, and Michael Bowling. The hanabi challenge: A new frontier for ai research. *Artificial Intelligence*, 280:103216, 2020. ISSN 0004-3702. doi: https://doi.org/10.1016/j.artint.2019.103216. URL https://www.sciencedirect.com/ science/article/pii/S0004370219300116.

- 546 Jakob Bauer, Kate Baumli, Feryal Behbahani, Avishkar Bhoopchand, Nathalie Bradley-Schmieg, 547 Michael Chang, Natalie Clay, Adrian Collister, Vibhavari Dasagi, Lucy Gonzalez, Karol Gre-548 gor, Edward Hughes, Sheleem Kashem, Maria Loks-Thompson, Hannah Openshaw, Jack Parker-549 Holder, Shreya Pathak, Nicolas Perez-Nieves, Nemanja Rakicevic, Tim Rocktäschel, Yannick 550 Schroecker, Satinder Singh, Jakub Sygnowski, Karl Tuyls, Sarah York, Alexander Zacherl, and 551 Lei M Zhang. Human-timescale adaptation in an open-ended task space. In Andreas Krause, 552 Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett 553 (eds.), Proceedings of the 40th International Conference on Machine Learning, volume 202 of 554 Proceedings of Machine Learning Research, pp. 1887–1935. PMLR, 23–29 Jul 2023. URL https://proceedings.mlr.press/v202/bauer23a.html. 555
- Michael Beukman, Samuel Coward, Michael Matthews, Mattie Fellows, Minqi Jiang, Michael Dennis, and Jakob N. Foerster. Refining minimax regret for unsupervised environment design. *CoRR*, abs/2402.12284, 2024. doi: 10.48550/ARXIV.2402.12284. URL https://doi.org/10.48550/arXiv.2402.12284.
- Matteo Bortoletto, Constantin Ruhdorfer, Adnen Abdessaied, Lei Shi, and Andreas Bulling. Limits
 of theory of mind modelling in dialogue-based collaborative plan acquisition. In *Proc. 62nd Annual Meeting of the Association for Computational Linguistics (ACL)*, pp. 1–16, 2024a.
 - Matteo Bortoletto, Lei Shi, and Andreas Bulling. Neural reasoning about agents' goals, preferences, and actions. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(1):456–464, Mar. 2024b. doi: 10.1609/aaai.v38i1.27800. URL https://ojs.aaai.org/index.php/AAAI/article/view/27800.
 - James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL http: //github.com/google/jax.
- Jake Bruce, Michael Dennis, Ashley Edwards, Jack Parker-Holder, Yuge Shi, Edward Hughes, Matthew Lai, Aditi Mavalankar, Richie Steigerwald, Chris Apps, Yusuf Aytar, Sarah Bechtle, Feryal Behbahani, Stephanie Chan, Nicolas Heess, Lucy Gonzalez, Simon Osindero, Sherjil Ozair, Scott Reed, Jingwei Zhang, Konrad Zolna, Jeff Clune, Nando de Freitas, Satinder Singh, and Tim Rocktäschel. Genie: Generative interactive environments, 2024. URL https://arxiv.org/abs/2402.15391.
 - Micah Carroll, Rohin Shah, Mark K Ho, Tom Griffiths, Sanjit Seshia, Pieter Abbeel, and Anca Dragan. On the utility of learning about humans for human-ai coordination. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), Advances in Neural Information Processing Systems, volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper/2019/file/ f5b1b89d98b7286673128a5fb112cb9a-Paper.pdf.
- Karl Cobbe, Oleg Klimov, Chris Hesse, Taehoon Kim, and John Schulman. Quantifying generalization in reinforcement learning. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 1282–1289. PMLR, 09–15 Jun 2019. URL https: //proceedings.mlr.press/v97/cobbe19a.html.
- Karl Cobbe, Christopher Hesse, Jacob Hilton, and John Schulman. Leveraging procedural generation to benchmark reinforcement learning. In *Proceedings of the 37th International Conference on Machine Learning*, ICML'20. JMLR.org, 2020.

- Allan Dafoe, Edward Hughes, Yoram Bachrach, Tantum Collins, Kevin R. McKee, Joel Z. Leibo, Kate Larson, and Thore Graepel. Open Problems in Cooperative AI, December 2020. URL http://arxiv.org/abs/2012.08630.
- Michael Dennis, Natasha Jaques, Eugene Vinitsky, Alexandre Bayen, Stuart Russell, Andrew Critch, and Sergey Levine. Emergent complexity and zero-shot transfer via unsupervised environment design. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS'20, Red Hook, NY, USA, 2020. Curran Associates Inc. ISBN 9781713829546.
- Lin Ding, Yong Tang, Tao Wang, Tianle Xie, Peihao Huang, and Bingsan Yang. A Cooperative Decision-Making Approach Based on a Soar Cognitive Architecture for Multi-Unmanned Vehicles. Drones, 8(4):155, April 2024. ISSN 2504-446X. doi: 10.3390/drones8040155. URL https://www.mdpi.com/2504-446X/8/4/155.
- Jesse Farebrother, Marlos C. Machado, and Michael H. Bowling. Generalization and regularization in dqn. ArXiv, abs/1810.00123, 2018. URL https://api.semanticscholar.org/ CorpusID:52904113.
- Jakob Foerster, Ioannis Alexandros Assael, Nando de Freitas, and Shimon Whiteson. Learning to communicate with deep multi-agent reinforcement learning. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett (eds.), Advances in Neural Information Processing Systems, volume 29.
 Curran Associates, Inc., 2016. URL https://proceedings.neurips.cc/paper_files/paper/2016/file/c7635bfd99248a2cdef8249ef7bfbef4-Paper.pdf.
- Jakob N. Foerster, Gregory Farquhar, Triantafyllos Afouras, Nantas Nardelli, and Shimon Whiteson. Counterfactual multi-agent policy gradients. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence and Thirtieth Innovative Applications of Artificial Intelligence Conference and Eighth AAAI Symposium on Educational Advances in Artificial Intelligence*, AAAI'18/IAAI'18/EAAI'18. AAAI Press, 2018. ISBN 978-1-57735-800-8.
- Matthew Fontaine, Ya-Chuan Hsu, Yulun Zhang, Bryon Tjanaka, and Stefanos Nikolaidis. On the Importance of Environments in Human-Robot Coordination. In *Proceedings of Robotics: Science and Systems*, Virtual, July 2021. doi: 10.15607/RSS.2021.XVII.038.
- Kunihiko Fukushima. Cognitron: A self-organizing multilayered neural network. *Biological Cybernetics*, 20:121–136, 1975. URL https://api.semanticscholar.org/CorpusID: 28586460.
- Kanishk Gandhi, Gala Stojnic, Brenden M. Lake, and Moira Dillon. Baby intuitions benchmark
 (BIB): Discerning the goals, preferences, and actions of others. In *Thirty-Fifth Conference on Neural Information Processing Systems*, 2021. URL https://arxiv.org/abs/2102.
 11938.
- Alex Graves, Marc G. Bellemare, Jacob Menick, Rémi Munos, and Koray Kavukcuoglu. Automated curriculum learning for neural networks. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, ICML'17, pp. 1311–1320. JMLR.org, 2017.
- Albert Gu, Karan Goel, and Christopher Re. Efficiently modeling long sequences with structured state spaces. In International Conference on Learning Representations, 2022. URL https://openreview.net/forum?id=uYLFoz1vlAC.
- Jonathan Heek, Anselm Levskaya, Avital Oliver, Marvin Ritter, Bertrand Rondepierre, Andreas
 Steiner, and Marc van Zee. Flax: A neural network library and ecosystem for JAX, 2023. URL
 http://github.com/google/flax.
- Johannes Heinrich, Marc Lanctot, and David Silver. Fictitious self-play in extensive-form games. In Francis Bach and David Blei (eds.), *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pp. 805–813, Lille, France, 07–09 Jul 2015. PMLR. URL https://proceedings.mlr.press/v37/ heinrich15.html.
- 647 Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9:1735–80, 12 1997. doi: 10.1162/neco.1997.9.8.1735.

686

687

688

689

648	Hengyuan Hu, Adam Lerer, Alex Peysakhovich, and Jakob Foerster, "Other-play" for zero-shot
649	coordination In Hal Daumé III and Aarti Singh (eds.) Proceedings of the 37th International
650	Conference on Machine Learning, volume 119 of Proceedings of Machine Learning Research.
651	pp. 4399–4410. PMLR. 13–18 Jul 2020.
652	FF. 1000 Contract, of Contraction

- Nick Jakobi. Evolutionary robotics and the radical envelope-of-noise hypothesis. *Adaptive Behavior*, 6(2):325–368, September 1997. ISSN 1741-2633. doi: 10.1177/105971239700600205. URL http://dx.doi.org/10.1177/105971239700600205.
- Minqi Jiang, Michael D Dennis, Jack Parker-Holder, Jakob Nicolaus Foerster, Edward Grefen stette, and Tim Rocktäschel. Replay-guided adversarial environment design. In A. Beygelzimer,
 Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.), Advances in Neural Information Processing
 Systems, 2021a. URL https://openreview.net/forum?id=5UZ-AcwFDKJ.
- Minqi Jiang, Edward Grefenstette, and Tim Rocktäschel. Prioritized level replay. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 4940–4950. PMLR, 18–24 Jul 2021b. URL https://proceedings.mlr.press/v139/jiang21b.html.
- Minqi Jiang, Michael Dennis, Edward Grefenstette, and Tim Rocktäschel. minimax: Efficient base lines for autocurricula in JAX. *CoRR*, abs/2311.12716, 2023. doi: 10.48550/ARXIV.2311.12716.
 URL https://doi.org/10.48550/arXiv.2311.12716.
- Rawal Khirodkar, Donghyun Yoo, and Kris M. Kitani. Adversarial domain randomization. CoRR, abs/1812.00491v2, 2018. URL https://arxiv.org/abs/1812.00491v2.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua Bengio and Yann LeCun (eds.), 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings, 2015. URL http://arxiv.org/abs/1412.6980.
- Marc Lanctot, Vinicius Zambaldi, Audrunas Gruslys, Angeliki Lazaridou, Karl Tuyls, Julien Pero lat, David Silver, and Thore Graepel. A unified game-theoretic approach to multiagent reinforce ment learning. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan,
 and R. Garnett (eds.), Advances in Neural Information Processing Systems, volume 30. Cur ran Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper_files/
 paper/2017/file/3323fel1e9595c09af38fe67567a9394-Paper.pdf.
- Adam Lerer and Alexander Peysakhovich. Learning existing social conventions via observationally augmented self-play. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*, AIES '19, pp. 107–114, New York, NY, USA, 2019. Association for Computing Machinery. ISBN 9781450363242. doi: 10.1145/3306618.3314268. URL https://doi.org/10. 1145/3306618.3314268.
 - Wenjun Li, Pradeep Varakantham, and Dexun Li. Generalization through diversity: improving unsupervised environment design. In *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence*, IJCAI '23, 2023a. ISBN 978-1-956792-03-4. doi: 10.24963/ijcai.2023/601. URL https://doi.org/10.24963/ijcai.2023/601.
- Yang Li, Shao Zhang, Jichen Sun, Yali Du, Ying Wen, Xinbing Wang, and Wei Pan. Cooperative open-ended learning framework for zero-shot coordination. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 20470–20484. PMLR, 23–29 Jul 2023b. URL https://proceedings.mlr.press/v202/li23au.html.
- Jijia Liu, Chao Yu, Jiaxuan Gao, Yuqing Xie, Qingmin Liao, Yi Wu, and Yu Wang. Llm-powered hierarchical language agent for real-time human-ai coordination. In Mehdi Dastani, Jaime Simão Sichman, Natasha Alechina, and Virginia Dignum (eds.), *Proceedings of the 23rd International Conference on Autonomous Agents and Multiagent Systems, AAMAS 2024, Auckland, New Zealand, May 6-10, 2024*, pp. 1219–1228. ACM, 2024. doi: 10.5555/3635637.3662979. URL https://dl.acm.org/doi/10.5555/3635637.3662979.

702	Ryan Lowe, YI WU, Aviv Tamar, Jean Harb, OpenAI Pieter Abbeel, and Igor Mordatch. Multi-
703	agent actor-critic for mixed cooperative-competitive environments. In I. Guyon, U. Von
704	Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), Ad-
705	vances in Neural Information Processing Systems, volume 30. Curran Associates, Inc.,
706	2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/
707	file/68a9750337a418a86fe06c1991a1d64c-Paper.pdf.

- Chris Lu, Yannick Schroecker, Albert Gu, Emilio Parisotto, Jakob Foerster, Satinder Singh, and Feryal Behbahani. Structured state space models for in-context reinforcement learning. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), Advances in Neural Information Processing Systems, volume 36, pp. 47016–47031. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/ file/92d3d2a9801211ca3693ccb2faa1316f-Paper-Conference.pdf.
- Michael Matthews, Michael Beukman, Benjamin Ellis, Mikayel Samvelyan, Matthew Thomas Jackson, Samuel Coward, and Jakob N. Foerster. Craftax: A lightning-fast benchmark for open-ended reinforcement learning. *CoRR*, abs/2402.16801, 2024. doi: 10.48550/ARXIV.2402.16801. URL https://doi.org/10.48550/arXiv.2402.16801.
- Vinod Nair and Geoffrey E. Hinton. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th International Conference on International Conference on Machine Learning*, ICML'10, pp. 807–814, Madison, WI, USA, 2010. Omnipress. ISBN 9781605589077.
- Alex Nichol, Vicki Pfau, Christopher Hesse, Oleg Klimov, and John Schulman. Gotta learn fast: A new benchmark for generalization in RL. *CoRR*, abs/1804.03720, 2018. URL http://arxiv.org/abs/1804.03720.
- Stefanos Nikolaidis and Julie Shah. Human-robot cross-training: computational formulation, modeling and evaluation of a human team training strategy. In *Proceedings of the 8th ACM/IEEE International Conference on Human-Robot Interaction*, HRI '13, pp. 33–40. IEEE Press, 2013. ISBN 9781467330558.
- Alexander Nikulin, Vladislav Kurenkov, Ilya Zisman, Viacheslav Sinii, Artem Agarkov, and Sergey Kolesnikov. XLand-minigrid: Scalable meta-reinforcement learning environments in JAX. In Intrinsically-Motivated and Open-Ended Learning Workshop @NeurIPS2023, 2023. URL https://openreview.net/forum?id=xALDC4aHGz.
- Johan S. Obando-Ceron, Ghada Sokar, Timon Willi, Clare Lyle, Jesse Farebrother, Jakob N. Foerster, Gintare Karolina Dziugaite, Doina Precup, and Pablo Samuel Castro. Mixtures of experts unlock parameter scaling for deep RL. *CoRR*, abs/2402.08609, 2024. doi: 10.48550/ARXIV. 2402.08609. URL https://doi.org/10.48550/arXiv.2402.08609.
 - Frans A. Oliehoek and Christopher Amato. A Concise Introduction to Decentralized POMDPs. Springer International Publishing, 2016. ISBN 9783319289298. doi: 10.1007/978-3-319-28929-8. URL http://dx.doi.org/10.1007/978-3-319-28929-8.

739

- Thomas O'neill, Nathan McNeese, Amy Barron, and Beau Schelble. Human-autonomy teaming: A review and analysis of the empirical literature. *Human Factors The Journal of the Human Factors and Ergonomics Society*, 64:1–35, 10 2020. doi: 10.1177/0018720820960865.
- Jack Parker-Holder, Minqi Jiang, Michael Dennis, Mikayel Samvelyan, Jakob Foerster, Edward Grefenstette, and Tim Rocktäschel. Evolving curricula with regret-based environment design. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 17473–17498. PMLR, 17–23 Jul 2022.
- David Premack and G. Woodruff. Does the chimpanzee have a theory of mind? *Behavioral and Brain Sciences*, 4(4):515–629, 1978. doi: 10.1017/s0140525x00076512.
- Neil Rabinowitz, Frank Perbet, Francis Song, Chiyuan Zhang, S. M. Ali Eslami, and Matthew
 Botvinick. Machine theory of mind. In Jennifer Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 4218–4227. PMLR, 10–15 Jul 2018. URL https://proceedings.
 mlr.press/v80/rabinowitz18a.html.

- David G. Rand and Martin A. Nowak. Human cooperation. *Trends in Cognitive Sciences*, 17(8):413–425, August 2013. ISSN 1364-6613, 1879-307X. doi: 10.1016/j.tics.2013.06.
 URL https://www.cell.com/trends/cognitive-sciences/abstract/ S1364-6613(13)00121-6.
- Constantin Ruhdorfer. Into the minds of the chefs: Using theory of mind for robust collaboration with humans in overcooked. Master's thesis, 2023. URL http://elib.uni-stuttgart. de/handle/11682/13983.
- Alexander Rutherford, Michael Beukman, Timon Willi, Bruno Lacerda, Nick Hawes, and Jakob
 Foerster. No regrets: Investigating and improving regret approximations for curriculum discovery,
 2024a. URL https://arxiv.org/abs/2408.15099.
- Alexander Rutherford, Benjamin Ellis, Matteo Gallici, Jonathan Cook, Andrei Lupu, Garar Ingvarsson, Timon Willi, Akbir Khan, Christian Schroeder de Witt, Alexandra Souly, Saptarashmi Bandyopadhyay, Mikayel Samvelyan, Minqi Jiang, Robert Lange, Shimon Whiteson, Bruno Lacerda, Nick Hawes, Tim Rocktäschel, Chris Lu, and Jakob Foerster. Jaxmarl: Multi-agent rl environments and algorithms in jax. In *Proceedings of the 23rd International Conference on Autonomous Agents and Multiagent Systems*, AAMAS '24, pp. 2444–2446, Richland, SC, 2024b. International Foundation for Autonomous Agents and Multiagent Systems. ISBN 9798400704864.
- Dorsa Sadigh, Shankar Sastry, Sanjit A. Seshia, and Anca D. Dragan. Planning for autonomous cars that leverage effects on human actions. In *Robotics: Science and Systems*, 2016. URL https://api.semanticscholar.org/CorpusID:7087988.
- Mikayel Samvelyan, Tabish Rashid, Christian Schroeder de Witt, Gregory Farquhar, Nantas Nardelli, Tim G. J. Rudner, Chia-Man Hung, Philip H. S. Torr, Jakob Foerster, and Shimon Whiteson. The starcraft multi-agent challenge. In *Proceedings of the 18th International Conference on Autonomous Agents and MultiAgent Systems*, AAMAS '19, pp. 2186–2188, Richland, SC, 2019. International Foundation for Autonomous Agents and Multiagent Systems. ISBN 9781450363099.
- Mikayel Samvelyan, Akbir Khan, Michael D Dennis, Minqi Jiang, Jack Parker-Holder, Jakob Nicolaus Foerster, Roberta Raileanu, and Tim Rocktäschel. MAESTRO: Open-ended environment design for multi-agent reinforcement learning. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id= sKWlRDzPfd7.
- Wilko Schwarting, Tim Seyde, Igor Gilitschenski, Lucas Liebenwein, Ryan M Sander, Sertac Karaman, and Daniela Rus. Deep latent competition: Learning to race using visual control policies in latent space. In *Conference on Robot Learning*, 2021. URL https://api.semanticscholar.org/CorpusID:231979320.

794

796

797

798

- Jimmy T.H. Smith, Andrew Warrington, and Scott Linderman. Simplified state space layers for sequence modeling. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=Ai8Hw3AXqks.
- Peter Stone, Gal Kaminka, Sarit Kraus, and Jeffrey Rosenschein. Ad hoc autonomous agent teams: Collaboration without pre-coordination. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 24, pp. 1504–1509, 2010.
- DJ Strouse, Kevin McKee, Matt Botvinick, Edward Hughes, and Richard Everett. Collaborating with humans without human data. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), Advances in Neural Information Processing Systems, volume 34, pp. 14502–14515. Curran Associates, Inc., 2021. URL https://proceedings.neurips.cc/paper_files/paper/2021/file/797134c3e42371bb4979a462eb2f042a-Paper.pdf.
- Joseph Suarez, Yilun Du, Clare Zhu, Igor Mordatch, and Phillip Isola. The neural mmo platform for massively multiagent research. In J. Vanschoren and S. Yeung (eds.), *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*, volume 1, 2021. URL https://datasets-benchmarks-proceedings.neurips.cc/paper/2021/ file/44f683a84163b3523afe57c2e008bc8c-Paper-round1.pdf.

810 Joseph Suarez, Kyoung Whan Choe, David Bloomin, Hao Xiang Li, Nikhil Pinnaparaju, Nishaanth 811 Kanna, Daniel Scott, Ryan Sullivan, Rose Shuman, Lucas de Alcantara, Herbie Bradley, 812 Chenghui Yu, Yuhao Jiang, Qimai Li, Jiaxin Chen, Xiaolong Zhu, Louis Castricato, and Phillip 813 Isola. Neural mmo 2.0: A massively multi-task addition to massively multi-agent learning. 814 In A. Oh, T. Neumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), Advances in Neural Information Processing Systems, volume 36, pp. 50094-50104. Curran Associates, 815 Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/ 816 2023/file/9ca22870ae0ba55ee50ce3e2d269e5de-Paper-Datasets_and_ 817 Benchmarks.pdf. 818

- Weihao Tan, Wentao Zhang, Shanqi Liu, Longtao Zheng, Xinrun Wang, and Bo An. True knowledge comes from practice: Aligning large language models with embodied environments via reinforce-ment learning. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=hILVmJ4Uvu.
- Open Ended Learning Team, Adam Stooke, Anuj Mahajan, Catarina Barros, Charlie Deck, Jakob
 Bauer, Jakub Sygnowski, Maja Trebacz, Max Jaderberg, Michaël Mathieu, Nat McAleese,
 Nathalie Bradley-Schmieg, Nathaniel Wong, Nicolas Porcel, Roberta Raileanu, Steph Hughes Fitt, Valentin Dalibard, and Wojciech Marian Czarnecki. Open-ended learning leads to generally
 capable agents, 2021. URL https://arxiv.org/abs/2107.12808.
- 829 Oriol Vinyals, Igor Babuschkin, Wojciech M. Czarnecki, Michaël Mathieu, Andrew Dudzik, Jun-830 young Chung, David Choi, Richard Powell, Timo Ewalds, Petko Georgiev, Junhyuk Oh, Dan 831 Horgan, Manuel Kroiss, Ivo Danihelka, Aja Huang, L. Sifre, Trevor Cai, John P. Agapiou, Max 832 Jaderberg, Alexander Sasha Vezhnevets, Rémi Leblond, Tobias Pohlen, Valentin Dalibard, David 833 Budden, Yury Sulsky, James Molloy, Tom Le Paine, Caglar Gulcehre, Ziyun Wang, Tobias Pfaff, Yuhuai Wu, Roman Ring, Dani Yogatama, Dario Wünsch, Katrina McKinney, Oliver Smith, Tom 834 Schaul, Timothy P. Lillicrap, Koray Kavukcuoglu, Demis Hassabis, Chris Apps, and David Sil-835 ver. Grandmaster level in starcraft ii using multi-agent reinforcement learning. Nature, 575:350 836 -354,2019. URL https://api.semanticscholar.org/CorpusID:204972004. 837
- Liza Vizmathy, Katarina Begus, Gunther Knoblich, György Gergely, and Arianna Curioni. Better Together: 14-Month-Old Infants Expect Agents to Cooperate. *Open Mind*, 8:1–16, February 2024. ISSN 2470-2986. doi: 10.1162/opmi_a_00115. URL https://doi.org/10.1162/ opmi_a_00115.
- Rui Wang, Joel Lehman, Jeff Clune, and Kenneth O. Stanley. Paired open-ended trailblazer (POET):
 endlessly generating increasingly complex and diverse learning environments and their solutions.
 CoRR, abs/1901.01753, 2019. URL http://arxiv.org/abs/1901.01753.
- Rui Wang, Joel Lehman, Aditya Rawal, Jiale Zhi, Yulun Li, Jeff Clune, and Kenneth O. Stanley. Enhanced poet: Open-ended reinforcement learning through unbounded invention of learning challenges and their solutions. In *International Conference on Machine Learning*, 2020. URL https://api.semanticscholar.org/CorpusID:213175678.
- Zifan Wu, Chao Yu, Deheng Ye, Junge Zhang, haiyin piao, and Hankz Hankui Zhuo. Coordinated proximal policy optimization. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.), Advances in Neural Information Processing Systems, 2021. URL https://openreview.net/forum?id=iCJFwoylT-q.
- Xue Yan, Jiaxian Guo, Xingzhou Lou, Jun Wang, Haifeng Zhang, and Yali Du. An efficient end-to-end training approach for zero-shot human-ai coordination. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), Advances in Neural Information Processing Systems, volume 36, pp. 2636–2658. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/07a363fd2263091c2063998e0034999c-Paper-Conference.pdf.
- Mesut Yang, Micah Carroll, and Anca D. Dragan. Optimal behavior prior: Data-efficient human models for improved human-ai collaboration. *CoRR*, abs/2211.01602, 2022. doi: 10.48550/ARXIV.2211.01602. URL https://doi.org/10.48550/arXiv.2211.01602.

- Chang Ye, Ahmed Khalifa, Philip Bontrager, and Julian Togelius. Rotation, translation, and cropping for zero-shot generalization. In *2020 IEEE Conference on Games (CoG)*, pp. 57–64, 2020. doi: 10.1109/CoG47356.2020.9231907.
- Chao Yu, Akash Velu, Eugene Vinitsky, Jiaxuan Gao, Yu Wang, Alexandre Bayen, and YI WU. The surprising effectiveness of ppo in cooperative multi-agent games. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), Advances in Neural Information Processing Systems, volume 35, pp. 24611–24624. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/ 9c1535a02f0ce079433344e14d910597-Paper-Datasets_and_Benchmarks. pdf.
- Chao Yu, Jiaxuan Gao, Weilin Liu, Botian Xu, Hao Tang, Jiaqi Yang, Yu Wang, and Yi Wu. Learning zero-shot cooperation with humans, assuming humans are biased. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum? id=TrwE819aJzs.
- Amy Zhang, Nicolas Ballas, and Joelle Pineau. A dissection of overfitting and generalization in continuous reinforcement learning. *CoRR*, abs/1806.07937, 2018a. URL http://arxiv.org/abs/1806.07937.
- Chiyuan Zhang, Oriol Vinyals, Rémi Munos, and Samy Bengio. A study on overfitting in deep reinforcement learning. *CoRR*, abs/1804.06893, 2018b. URL http://arxiv.org/abs/1804.06893.
 - Rui Zhao, Jinming Song, Yufeng Yuan, Haifeng Hu, Yang Gao, Yi Wu, Zhongqian Sun, and Wei Yang. Maximum entropy population-based training for zero-shot human-ai coordination. *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(5):6145–6153, Jun. 2023. doi: 10. 1609/aaai.v37i5.25758. URL https://ojs.aaai.org/index.php/AAAI/article/ view/25758.

A APPENDIX

918 A.1 ACCESSIBILITY OF THE BENCHMARK

920 We make our challenge available under the Apache License 2.0 via a code repository: https: 921 //anonymised.edu. Our environment is built on top of the existing minimax project (accessible under Apache License 2.0 via https://github.com/facebookresearch/minimax) 922 and is thus accessible to researchers who are already familiar with the project. minimax is exten-923 sively documented, fast, and supports multi-device training. For all details, including a full descrip-924 tion of the advantages of minimax, we kindly refer the reader to the accompanying publication 925 (Jiang et al., 2023). Our Overcooked adaption is extended from the one in JaxMARL also acces-926 sible under Apache License 2.0 via https://github.com/FLAIROx/JaxMARL. Our code 927 includes extensive documentation and examples of how it may be used. Additionally, our code is 928 written in a modular fashion and other multi-agent environments can be integrated with the runners 929 thanks to the careful design of the original project. 930

A.2 BROADER IMPACTS

933 While our work is largely foundational and concerned with providing the research community with 934 the appropriate tools for the training and evaluation of agents in game-like environments, special 935 caution is always imperative should this research be applied to human-AI collaboration. Even though 936 our goal is to improve collaboration, safeguards should be applied to make sure that humans are always safe from harm. Especially so in real-world applications where accidents could potentially 937 result in bodily harm. Since our work is still far removed from any real-world application, we do 938 not expect that our work in its present form carries the risk of materialising these harms. Some form 939 of unsupervised environment design in collaborative environments might be part of future systems 940 and we therefore acknowledge these risks. This work of course also carries the potential to improve 941 human-AI collaboration and we make an important contribution to advancing the field with potential 942 impacts in all kinds of human-machine interaction.

943 944

931

932

945 A.3 INFRASTRUCTURE & TOOLS

946 We ran our experiments on a server running Ubuntu 22.04, equipped with NVIDIA Tesla V100-947 SXM2 GPUs with 32GB of memory and Intel Xeon Platinum 8260 CPUs. All training runs are exe-948 cuted on a single GPU only. We trained our models using Jax (Bradbury et al., 2018) and Flax (Heek 949 et al., 2023) with 1, 2 and 3 as random seed for training DCD methods and 1 to 8 as random seeds for the populations. Training the DCD methods usually finishes in under 24 hours, only SoftMoE 950 and PAIRED-based methods take longer. SoftMoE-based policies often take an extra 50% wall-951 clock time to train. Noticeable is also that our S5 implementation is the fastest, usually needing 952 30% less time. Both are compared to the default architectures' training time. In the longest case, 953 the combination of a SoftMoE-LSTM policy trained with PAIRED takes about 80 hours to complete 954 training. Our benchmark should be runnable on any system that features a single CUDA-compatible 955 GPU. Although in our experience our experiments will require 32GB VRAM to run. 956

950 957 958

A.4 HYPERPARAMETERS

959 We overview all hyperparameters for training in Table 4 and provide details on the hyperparameter 960 search used in Table 5. This search was conducted on smaller single layout runs to determine 961 reasonable values as complete runs would have been computationally infeasible. Furthermore we 962 show the hyperparameters for each DCD method separately: DR hyperparameters in Table 6, PLR 963 hyperparameters in Table 7, ACCEL hyperparameters in Table 8, and PAIRED hyperparameters in 964 Table 9. DR hyperparameters govern how Overcooked levels are generated randomly and apply to 965 all other processes in which a random level is sampled, for instance, in PLR, in which case the same hyperparameters apply. 966

967

968 A.5 NEURAL NETWORK ARCHITECTURES

969

This work employs an actor-critic architecture using a separate actor and critic in which the critic is centralised for training via MAPPO (Yu et al., 2022). For the actor, the observations are of shape $h \times w \times 26$, while for the centralised critic, we concatenate the observations along the last axis to

973	• • •	
974	Description	Value
975	Ontimizor	Adam (Vingma & Pa 2015)
976		Adalli (Kiligilia & Ba, 2013)
977	Adam p_1	0.9
079	Adam β_2	0.999
510	Adam ϵ	$1 \cdot 10^{-5}$
979	Learning Rate η	$3 \cdot 10^{-4}$
980	Learning Rate Annealing	-
981	Max Grad Norm	0.5
982	Discount Pate of	0.000
983	CAE	0.999
984		0.98
005	Entropy Coefficient	0.01
985	Value Loss Coefficient	0.5
986	# PPO Epochs	8
987	# PPO Minibatches	4
988	# PPO Steps	400
989	PPO Value Loss	Clipped
990	PPO Value Loss Clip Value	0.2
991	Reward Shaping	Yes (linearly decreased over training)

Table 4: Hyperparamters of the learning process.

Table 5: Values used for a grid search over hyperparameters governing the learning process. Finally used values appear in **bold**.

995		
996	Description	Value
997	Learning Rate η	$[1 \cdot 10^{-4}, 3 \cdot \mathbf{10^{-4}}, 5 \cdot 10^{-4}, 1 \cdot 10^{-3}]$
998		[0.01.0.1]
999	Entropy Coefficient	[0.01 0.1]
1000	# PPO Steps	[256, 400]
1001	# Hidden Layers	[2, 3, 4]
1001	Reward Shaping Annealing Steps	[0, 2500000, 5000000, until end]
1002		

1003 1004

992 993

994

997

972

form a centralised observation, i.e. the centralised observation has shape $h \times w \times 52$ following prior 1005 work (Yu et al., 2023).

All our networks feature a convolutional encoder f_c . This encoder always features three 2D convo-1007 lutions of 32, 64 and 32 channels with kernel size 3×3 each and pads the input with zeros. Our 1008 default activation function is ReLU (Fukushima, 1975; Nair & Hinton, 2010) which we apply after 1009 every convolutional block. We feed the output of f_c to a feed-forward neural network f_e with three 1010 layers with 64 neurons, ReLU and LayerNorm (Ba et al., 2016) applied each. f_e takes the flattened 1011 representation produced by f_c and produces an embedding $e \in \mathbb{R}^{b \times t \times 64}$ that we feed into a recur-1012 rent neural network (either LSTM (Hochreiter & Schmidhuber, 1997) or S5 (Smith et al., 2023)) 1013 to aggregate information along the temporal axis. We use this resulting embedding $e_t \in \mathbb{R}^{b \times 64}$ to 1014 produce action logits $l \in \mathbb{R}^{b \times 6}$ to parameterise a categorical distribution in the actor-network or directly produce a value $v \in \mathbb{R}^{b \times 1}$ in the critic network using a final projection layer. This archi-1015 1016 tecture is inspired by previous work on Overcooked-AI, specifically (Yu et al., 2023), see Figure 8 1017 for an overview. We also test the use of a S5 layer (Smith et al., 2023) in which case we use 2 S5 blocks, 2 S5 layers, use LayerNorm before the SSM block and the activation function described in 1018 the original work, i.e. $a(x) = \text{GELU}(x) \odot \sigma(W * \text{GELU}(x))$. 1019

1020 In the case of the SoftMoE architecture, we follow the same approach as in (Obando-Ceron et al., 1021 2024) and replace the penultimate layer with a SoftMoE layer. As in their work we use the PerConv tokenisation technique, i.e. given input $x \in \mathbb{N}^{h \times w \times 26}$ we take the output $y \in \mathbb{R}^{h \times w \times 32}$ of f_c and 1023 construct $h \times w$ tokens with dimension d = 32 that we then feed into the SoftMoE layer. We always use 32 slots and 4 experts for this layer, see (Obando-Ceron et al., 2024) for details on this layer. The 1024 resulting embedding is then passed into the two remaining linear layers before being also passed to 1025 RNN and used to produce an action or value, equivalent to the description above, compare Figure 9.

1026			Table 6:	DR hyperparameters.	
1027					_
1028			Description	Value	
1029			n walls to place	Sampled uniformly between $0 - 15$	_
1030			n onion piles to place	Sampled uniformly between $1-2$	
1031			n plate piles to place	Sampled uniformly between $1-2$	
1032			<i>n</i> pots to place	Sampled uniformly between $1-2$	
1033			n goals to place	Sampled uniformly between $1-2$	
1034					-
1035		Table 7.	PIP specific hyperpara	maters in addition to the DR hyperna	romotors
1036			r LK specific hyperpara	meters in addition to the DK hyperpa	Tameters.
1037			Description	Valua	
1038			Description	value	
1039			UED Score	MaxMC (Jiang et al., 2021a)	
1040			PLR replay probabilit	$p \rho = 0.5$	
1041			PLR buffer size	4,000	
1042			PLR staleness coeffici	lent 0.3	
1043			PLR temperature PLR score ranks	0.1 Ves	
1044			PLR minimum fill rat	io 0.5	
1045			PI \mathbb{R}^{\perp}	Ves	
1046				Ves	
1047			PLR force unique leve	el Yes	
1048					
1049					
1050	Loctly	wa dagamih	a ann naturalia in tamac a	of nonemator count in Table 10	
1051	Lastry,	we describe		n parameter count in Table 10.	
1052					
1053	A.6	ADDITION	AL ANALYSIS		
1054		_		_	
1055	A.6.1	EVIDENC	E OF OVERFITTING IN C	JVERCOOKED AGENTS	
1056	To veri	ify that are	ents indeed overfit their	training layout in Overcooked we t	present Table 11 in
1057	which	we experim	ent with our weakest per	forming policy architecture, the CN	N-LSTM. This is to
1058	be expe	ected but ve	rifying is nonetheless im	portant.	
1059	1		, 6	1	
1060	162	DEDEODA	AANCE ACDOSS LEVELS		
1061	A.0.2	I EKFUKN	TANCE ACROSS LEVELS		
1062	To acco	ompany the	e overall performance me	easured by reward in the main paper	in Table 3 we also
1063	measur	e the mean	solved rate on display it	in Table 12.	
1064			1 .		
1065	163	DEDEODA	AANCE ON INDIVIDUAL	LEVELS	
1066	A.0.5	I ERFORM	TANCE ON INDIVIDUAL		
1067	We list	the perforn	nance of every individual	method on every single layout in Tab	le 13. Most notable
1068	is that	some layou	its are harder to learn th	an others. Our agents especially see	em to struggle with
1069	layouts	requiring	more complex forms of i	interaction, i.e. Coordination Ring, G	Counter Circuit and
1070	Forced	Coordinati	on. Forced Coordination	n especially seems difficult to solve	as no run achieves
1071	noticea	ble perforn	nance on it. This might	be due to the specific features of the	ayout, i.e. agents
1072	have ac	ccess to sev	eral objects and need to h	and them over the counter to produce	e any result.
1073					
1074	A.6.4	POPULAT	TION TRAINING DETAILS	3	
1075	m 1 1	· C . 1		, , , , , , , , , , , , , , , , , , ,	.1
1076	To both	n verity tha	t our implementation is	correct and to give an intuition into	the performance of
10//	une me	tion in E	e population, we present	t the training curves over all 8 seeds	or training an FCP

population in Figure 10. MEP was trained with exactly the same architecture and with the same amount of experience per agent. As in prior work (Zhao et al., 2023) we set the population entropy coefficient during training to $\alpha = 0.01$.

080	Table 8	8: ACCEL hyperparameters	in addition to the DR hyperparame	eters.
081				
082		Description	Value	
003		UED Score	MaxMC (Jiang et al., 2021a)	
004		PLR replay probability ρ	0.8	
000		PLR buffer size	4,000	
000		PLR staleness coefficient	0.3	
087		PLR temperature	0.1	
880		PLR score ranks	Yes	
089		PLR minimum fill ratio	0.5	
090			Yes	
091		PLR ^{II}	Yes	
092		PLR force unique level	Yes	
093		ACCEL Mutation	Overcooked Mutator	
094		ACCEL n initiations	20	
095		ACCEL subsample size	4	
096				
097	Table 9: PAIRED h	yperparameters. All PPO hy	yperparameters are the same betwee	in the student and
098	the teacher. The mi	inimax implementation for	llows to original one in (Dennis et	al., 2020) and we
099	SHCK TO IT TOO.			
100		acconnection	Valua	_
101	De	escription	value	
102	n :	students	2	
103	UI	ED Score I	Relative regret (Dennis et al., 2020)	
104	UI	ED first wall sets budget	Yes	
105	UI	ED noise dim	50	
106	PA	ARED Creator	OvercookedUED	
107				
109 110 111 112	A.6.5 DETAILED We present detailed through the average best on four of the f	RESULTS WITH POPULATI I zero-shot cooperation res of performance discussed in iva individual layouts in ter	IONS sults per layout in Tables 14 and 1 in the main text, we also find that P	15. As indicated AIRED performs
113	best on four of the f	ive marviauar rayouts in ter	ins of zero-shot cooperation.	
14 15	A.7 Error Ana	LYSIS NUMBERS		
116 117 118 119 120 121	We hypothesise tha To test this we desi systemic failures, i. sion. While method analysis suggests, th on each of the 24 la	t agents may fail to genera gn testing layouts that rotat e. cases in which an agent d is generally differed in how he low standard deviations youts, see Table 16.	alise since observations are hard to be and mirror features of the environ loes well in one environment but no well they performed along the sam show that any given method perform	generalise from. nment to look for t its mirrored ver- e line as previous ms similarly well
122 123	A.7.1 TRAINING	CURVES AND EVALUATIO	Ν	
124 125 126 127 128 129 130	In Figures 11, 12 an as well as the five u in Figure 11, the r Figure 13. Interesti reach the highest tr generalisation gap s	d 13 we show the returns of inseen evaluation levels. The sults for the S5 one in Fig ngly, while (SoftMoE) PAII aining returns, instead it ac small.	Four agent during training in both se ne results for the SoftMoE architect gure 12 and the results for the CN RED performs the best in our evalu chieves the highest training return v	een training levels ure are displayed IN-LSTM one in ations it does not vhile keeping the
131	A.8 VALIDATING	THE IMPLEMENTATION		
33	As an open-source lall the baselines.	benchmark, we emphasise a We do so in two important	correct implementation of the benc ways. Firstly, we base our imple	hmark, including mentation on the



implementation of the minimax benchmark (Jiang et al., 2023), making sure that we use publicly available code for all unsupervised environment design algorithms. Secondly, we test the implementation and adaption of the Overcooked-AI environment by fixing the generated training layouts to a single layout during training. This allows us to train on the 5 classic Overcooked layouts using our implementation. Our implementation is capable of solving these layouts, see Figure 10. We do this in part to argue for the fact that our benchmark is hard to solve and this is not a function of poorly configured or wrongly implemented algorithms.





Figure 11: Returns in training and evaluation levels over the duration of training for our SoftMoE architecture.



Figure 12: Returns in training and evaluation levels over the duration of training for our S5 architecture.





Table 10: Number of trainable parameters in each model.

	CNN-LSTM	SoftMoE-LSTM	CNN-S5
Parameter Count	197,254	316,102	193,670

1423Table 11: Comparing the layout a CNN-LSTM policy was trained on versus on which it was being
evaluated. The policies heavily overfit the training layout. All policies we tested exhibit this prop-
erty.1425erty.

	Asymm	Cramped	Counter	Forced	Coord
Asymm	343.4	0.0	0.0	0.0	0.0
Cramped	1.6	185.6	0.0	0.0	0.0
Counter	0.0	0.0	128.0	0.0	0.0
Forced	0.0	0.2	0.0	141.2	0.0
Coord	0.0	0.0	0.0	0.0	144.6

Table 12: Mean episode solved rate for the different methods averaged over the respective testing layouts. The best result is shown in **bold**. We report aggregate statistics over three random seeds. As a baseline we include an Oracle version for all architectures which was trained on the five testing layouts directly

Method	CNN-LSTM	SoftMoE-LSTM	CNN-S5
DR	$0.02\pm0.0\%$	$6.31\pm10.1\%$	$0.00\pm0.0\%$
$PLR^{\perp,\parallel}$	$0.00\pm0.0\%$	$0.33\pm0.3\%$	$0.00\pm0.0\%$
Pop. PAIRED	$0.00\pm0.0\%$	$14.62 \pm 7.6\%$	$0.00\pm0.0\%$
ACCEL∥	$0.00\pm0.0\%$	$0.08\pm0.1\%$	$0.00\pm0.0\%$
Oracle	$95.40\pm7.5\%$	$99.67\pm0.6\%$	$97.53 \pm 4.1\%$

1469Table 13: Performance on all evaluation layouts. We show the mean episode reward **R** and the mean1470episode solved rate **SR**. The overall best result per layout is presented in **bold** excluding oracle1471results.

Layout	Method	CNN-	LSTM	SoftMo	E-LSTM	CNN	N-S5
		R	SR	R	SR	R	SF
	DR	1.70	0.0%	1.54	0.2%	0.00	0.0%
	$\mathrm{PLR}^{\perp,\parallel}$	1.12	0.0%	5.02	2.1%	0.14	0.0%
Cramped	Pop. PAIRED	1.44	0.0%	37.02	57.7 %	0.50	0.0%
	ACCEL [∥]	0.92	0.0%	0.60	0.0%	0.60	0.0%
	Oracle	241.27	96.7%	245.54	100.0%	189.47	99.79
	DR	0.00	0.0%	0.00	0.0%	0.00	0.0%
	$PLR^{\perp,\parallel}$	0.00	0.0%	0.00	0.0%	0.00	0.0%
Coord	Pop. PAIRED	0.00	0.0%	16.78	14.6%	0.00	0.09
	ACCEL∥	0.00	0.0%	0.04	0.0%	0.02	0.09
	Oracle	197.8	100.0%	204.53	100.0%	119.33	99.09
	DR	0.00	0.0%	0.02	0.0%	0.00	0.09
	$PLR^{\perp,\parallel}$	0.00	0.0%	0.02	0.0%	0.02	0.04
Forced	Pop. PAIRED	0.00	0.0%	0.00	0.0%	0.00	0.09
	ACCEL∥	0.00	0.0%	0.00	0.0%	0.00	0.09
	Oracle	196.8	100.0%	204.53	100.0%	133.47	94.79
	DR	0.58	0.1%	8.64	4.4%	0.00	0.0%
	$PLR^{\perp,\parallel}$	0.08	0.0%	0.10	0.0%	0.08	0.0%
Asymm	Pop. PAIRED	0.28	0.0%	15.64	14.2%	0.08	0.09
	$\mathrm{ACCEL}^{\parallel}$	0.14	0.0%	0.04	0.0%	0.02	0.09
	Oracle	220.4	100.0%	277.8	98.4%	247.87	99.79
	DR	0.00	0.0%	0.00	0.0%	0.00	0.06
	$PLR^{\perp,\parallel}$	0.00	0.0%	0.00	0.0%	0.00	0.09
Counter	Pop. PAIRED	0.00	0.0%	1.38	0.0%	0.00	0.09
	ACCEL∥	0.00	0.0%	0.00	0.0%	0.00	0.00
	Oracle	91.2	77.3%	152.73	100.0%	84.93	94.79

Table 14: Zero-shot results using SoftMoE-LSTM policies playing with an FCP and MEP popula-tion of experts trained on the respective layout exclusively. We report the mean episode reward and standard deviation. The best result per layout is put in **bold**.

Method	Asymm	Counter	Cramped	Forced	Coord		
FCP							
Random	7.43 ± 12.19	8.89 ± 4.65	66.02 ± 38.28	1.95 ± 1.92	$20.49\pm$		
Stay	5.32 ± 12.07	0.38 ± 1.11	20.67 ± 33.05	0.00 ± 0.00	$0.95 \pm$		
Oracle	126.44 ± 27.13	22.63 ± 7.82	120.9 ± 10.86	22.08 ± 12.89	59.64 ± 2		
DR	18.18 ± 1.69	6.86 ± 5.27	65.05 ± 5.15	1.09 ± 0.21	17.88 ± 1		
$PLR^{\perp,\parallel}$	7.64 ± 0.89	5.60 ± 1.29	60.35 ± 6.89	1.76 ± 0.86	$21.90 \pm$		
Pop. PAIRED	24.51 ± 3.44	11.11 ± 1.67	81.92 ± 6.33	1.59 ± 0.57	29.72 ± 4		
ACCEL∥	8.60 ± 0.98	10.23 ± 0.85	65.46 ± 4.62	1.81 ± 1.25	$19.19\pm$		
		MI	EP				
Random	8.0 ± 9.12	22.46 ± 13.34	58.33 ± 34.83	2.55 ± 2.76	31.85 ± 1		
Stay	4.86 ± 7.21	5.2 ± 10.85	31.55 ± 47.13	0.0 ± 0.0	$1.53 \pm$		
Oracle	135.07 ± 30.27	39.33 ± 13.53	138.07 ± 10.0	56.1 ± 25.41	67.86 ± 1		
DR	19.32 ± 0.39	18.04 ± 5.75	62.77 ± 7.22	1.69 ± 0.67	$30.35 \pm$		
$PLR^{\perp,\parallel}$	7.53 ± 0.92	21.23 ± 1.91	57.2 ± 4.4	2.45 ± 1.23	$2.45 \pm$		
Pop. PAIRED	24.33 ± 2.27	23.72 ± 4.0	82.23 ± 9.38	2.96 ± 1.56	37.1 ± 0		
ACCEL∥	9.3 ± 0.71	18.33 ± 1.96	56.72 ± 4.15	2.21 ± 1.57	$28.52 \pm$		

Table 15: Zero-shot results using SoftMoE-LSTM policies playing with an FCP and MEP popu-lation of experts trained on the respective layout exclusively. We report the mean solved rate and standard deviation. The best result per layout is put in **bold**.

Method	Asymm	Counter	Cramped	Forced	Coord
Random	$8.52 \pm 17.52\%$	$5.00 \pm 6.70\%$	$69.43 \pm 38.45\%$	$0.00 \pm 0.00\%$	$30.89 \pm 3.$
Stay	$6.81 \pm 18.04\%$	$0.02 \pm 0.14\%$	$21.75 \pm 33.71\%$	$0.00\pm0.00\%$	0.14 ± 0.14
Oracle	$69.67 \pm 16.39\%$	$27.39 \pm 19.02\%$	$31.30 \pm 20.97\%$	$92.02 \pm 1.19\%$	96.96 ± 2
DR	$24.19 \pm 4.60\%$	$4.56\pm5.32\%$	$72.11 \pm 6.29\%$	$0.01 \pm 0.01\%$	23.76 ± 18
$PLR^{\perp,\parallel}$	$8.84 \pm 1.31\%$	$2.04\pm0.95\%$	$68.14 \pm 1.21\%$	$0.11 \pm 0.12\%$	30.89 ± 3
Pop. PAIRED	$32.48 \pm 4.00\%$	$7.91 \pm \mathbf{1.38\%}$	$85.54 \pm 6.08\%$	$0.09\pm0.07\%$	$f 48.31 \pm 11.$
ACCEL∥	$9.58\pm1.12\%$	$6.79 \pm 0.91\%$	$69.01 \pm 2.03\%$	$0.06\pm0.06\%$	24.13 ± 6
]	MEP		
Random	$9.25 \pm 2.02\%$	$36.04 \pm 4.38\%$	$67.75 \pm 5.48\%$	$0.00 \pm 0.00\%$	54.9 ± 5
Stay	$4.91\pm1.46\%$	$5.85 \pm 2.71\%$	$29.56 \pm 5.92\%$	$0.00\pm0.00\%$	1.02 ± 0
Oracle	$91.02 \pm 1.12\%$	$52.60 \pm 11.37\%$	$96.86 \pm 2.27\%$	$56.16 \pm 21.85\%$	75.23 ± 0
DR	$26.34 \pm 3.55\%$	$27.41 \pm 10.31\%$	$70.78 \pm 4.23\%$	$0.05 \pm 0.07\%$	50.07 ± 6
$PLR^{\perp,\parallel}$	$8.24 \pm 1.28\%$	$33.76 \pm 4.89\%$	$65.38 \pm 4.55\%$	$0.28\pm0.41\%$	50.97 ± 4
Pop. PAIRED	$32.79 \pm \mathbf{1.81\%}$	$36.48 \pm 8.14\%$	$80.60 \pm 6.74\%$	$0.38 \pm 0.51\%$	55.23 ± 8
ACCEL	$10.10 \pm 0.11\%$	$25.94 \pm 4.31\%$	$66.52 \pm 3.34\%$	$0.18 \pm 0.16\%$	48.80 ± 2

Method	SoftMoE-LSTM	CNN-S5	CNN-LSTM
DR	11.91 ± 0.8	0.00 ± 0.0	1.96 ± 0.3
$\mathrm{PLR}^{\perp,\parallel}$	4.39 ± 0.4	0.49 ± 0.2	0.65 ± 0.2
Pop. PAIRED	36.83 ± 1.1	0.58 ± 0.2	0.59 ± 0.1
ACCEL∥	2.91 ± 0.4	1.69 ± 0.3	0.78 ± 0.2

Table 16: Performance on mirrored and rotated levels, illustrated in Figure 6.