

# 000 CAYLEY MAZE: UNIVERSAL OPEN-ENDED REIN- 001 002 FORCEMENT LEARNING ENVIRONMENT 003 004

005 **Anonymous authors**

006 Paper under double-blind review

## 007 008 ABSTRACT 009

010  
011 Parametrizable environments with variable complexity are crucial for advancing  
012 fields such as Unsupervised Environment Design (UED), Open-Ended Learning,  
013 Curriculum Learning, and Meta Reinforcement Learning. However, the selection  
014 of environments in evaluation procedures, along with their complexities, is often  
015 either neglected or lacks formal justification. We propose the formal definition of  
016 complexity for Markov Decision Processes using Deterministic Finite Automata  
017 and Group Theory machinery. We introduce Cayley Maze, a novel open-ended  
018 reinforcement learning environment that naturally generalizes problems like solv-  
019 ing Rubik’s Cube or sorting. Cayley Maze is universal: every finite deterministic  
020 sparse MDP is an MDP of a certain instance of Cayley Maze. We demonstrate how  
021 Cayley Maze enables control over complexity, simplification, and combination of  
022 its instances.

## 023 024 1 INTRODUCTION

025 Designing agents capable of generalizing across diverse tasks and environments is both a challeng-  
026 ing and exciting problem in modern reinforcement learning. Open-ended learning, unsupervised  
027 environment design (UED), and curriculum learning remain attractive approaches to reach this goal  
028 (Hughes et al., 2024)

029 There has been a lot of progress in designing and implementing Open-Ended environments. For  
030 example, Genie (Bruce et al., 2024) or Craftax (Matthews et al., 2024). While the Genie is a huge  
031 achievement, we doubt that such an environment allows for producing challenging, algorithmic,  
032 or intelligent problems. On the other hand, some of the Unsupervised Environment Design(UED)  
033 algorithms(Beukman et al., 2024), (Parker-Holder et al., 2023) are still being evaluated on Minigrid  
034 (Boisvert et al., 2018) environments. In such an experiment-oriented field like Machine Learning,  
035 algorithms cannot be better than their evaluation procedures. We pose the question: what is a good  
036 evaluation procedure for Open Ended Learning? The partial answer is that variable complexity is  
037 necessarily its component. The core assumption of curriculum learning is the ability to produce  
038 observations that gradually become more complex. The diversity of parametrizable environments is  
039 an assumption of UED. Hence, it also relies on some notion of similarity/metric. This paper proposes  
040 a formal definition of the complexity of Reinforcement Learning environments. We discuss why the  
041 current heuristics, like state or action space cardinality, might be a naive estimate for this goal. We  
042 do it by translating certain concepts from the Language of Algebra, Deterministic Finite Automata  
043 theory, and Topology, and hopefully prove its usefulness. More than that, we introduce a novel  
044 Reinforcement Learning environment called Cayley Maze, which is, in a certain way, universal.  
045 We show that many important problems, such as sorting or Rubik’s Cube, are specific instances of  
046 Cayley Maze.

## 047 2 PRELIMINARIES

048 Here, we briefly recap some definitions of the algebraic approach to Automata theory. For more  
049 context we advise to read Pin (2021). *Monoid* is a set  $M$  equipped with the operation  $(\cdot) : M \times M \rightarrow M$   
050 for which the following hold:

- 054     • Identity element: there exists  $e \in M$  such that for every  $m \in M$   $e \cdot m = m \cdot e = m$ . The  
 055       identity element is always unique.  
 056     • Associativity: for every elements  $f, g, h \in M$ ,  $(f \cdot g) \cdot h = f \cdot (g \cdot h)$

058 A *group* is a monoid whose every element has an inverse, i.e. for every  $g \in G$  there exists  $g^{-1} \in G$   
 059 such that  $g \cdot g^{-1} = g^{-1} \cdot g = e$ .

060 Given a group  $G$ , a *G-action* on the set  $S$ , is the map  $\rho : G \times S \rightarrow S$ , for which  $\forall s \in S \rho(e, s) = s$ ,  
 061 and action commutes with multiplication:  $\rho(g, \rho(h, s)) = \rho(g \cdot h, s)$  for all  $g, h \in G$  and  $s \in S$ . A  
 062 *homomorphism* between monoids  $M$  and  $N$  is a map between its sets  $\phi : M \rightarrow N$ , such that for  
 063 every  $a, b \in M$   $\phi(a \cdot b) = \phi(a) \cdot \phi(b)$ . An *isomorphism* is a bijective homomorphism.

064 By  $[n]$  we denote an  $n$ -element set  $[n] = \{1, 2, \dots, n\}$ . Given set  $M$  and  $S$ , the set  $M^S$  is a set of  
 065 all functions from  $S$  to  $M$ . A monoid of functions from  $n$ -element set to itself with the operation of  
 066 composition is denoted by  $End([n]) = [n]^{[n]} = \{f : f : [n] \rightarrow [n]\}$ . The following theorem is the  
 067 core idea of our environment:

068 **Theorem 1** (Cayley's theorem). Every finite monoid  $N$  is isomorphic to some submonoid (subset)  
 069 of  $End([n])$  for certain  $n \in \mathbb{N}$ .

070 A subset  $G \subseteq M$  of monoid  $M$  is called a set of *generators* if it generates  $M$ , i.e. for every  $m \in M$   
 071 there is a sequence  $g_1, \dots, g_n$  such that  $m = g_n \cdot g_{n-1} \cdot \dots \cdot g_1$ . A free monoid on the set  $A$  is a  
 072 set  $M = A^*$ , containing all finite sequences of elements of  $A$ , with the operation of concatenation,  
 073 namely for  $x = x_m x_{m-1} \dots x_1$ ,  $y = y_n y_{n-1} \dots y_1$ ,  $x \cdot y = x_m x_{m-1} \dots x_1 y_n y_{n-1} \dots y_1$ .  
 074 A *congruence* on monoid  $M$  is an equivalence relation  $\sim$  on a set  $M$ , which is compatible with its  
 075 operation: for every  $a, b, c, d \in M$ ,  $a \sim c$ ,  $b \sim d \implies (a \cdot b) \sim (c \cdot d)$ . Every congruence induces  
 076 monoid structure on the set  $M / \sim$  of equivalence classes on  $M$  and a canonical homomorphism:  
 077  $\pi : M \rightarrow M / \sim$ ,  $\pi(a \cdot b) = [a \cdot b]_\sim = [a]_\sim \cdot [b]_\sim$ .  $M / \sim$  is called a *quotient monoid* of  $M$ . Given  
 078 a monoid  $M$  and its subset  $L$ , *syntactic congruence* on  $M$  is defined as  $a \sim_L b$  if for all  $x, y \in M$   
 079  $xay \in L \iff xby \in L$ . The quotient of this equivalence relation is called a *syntactic monoid*.

080 We call *Markov Decision Process* a tuple  $(A, S, s_0, R, T)$  where  $A$  is the set of actions,  $S$  - the set  
 081 of states,  $s_0$  - initial state,  $R : A \times S \rightarrow \mathbb{R}$  - reward function, and  $T : A \rightarrow (S \rightarrow Pr(S))$  transition  
 082 function, assigning the transition kernel  $T(a) = T_a$  on  $S$  to every state, where  $Pr(S)$  stands for the  
 083 space of probability distributions on the set  $S$ . We treat  $T_a$  as a matrix of size  $|S| \times |S|$ , whose value  
 084 at the row  $s_2$  and column  $s_1$  is denoted by  $T_a(s_1)(s_2)$  meaning the probability of getting from  $s_1$  to  
 085  $s_2$  by the action  $a$ . If  $T_a(s)$  are 0 - 1 valued for all  $a \in A$ ,  $s \in S$ ,  $T_a$  becomes a function  $S \rightarrow S$ ,  
 086 and we write  $T_a(s_1) = s_2$  instead of  $T_a(s_1)(s_2) = 1$ .

087 For the mathematical convenience, we assume that the action set  $A$  necessarily contains neutral (do  
 088 nothing) action  $e$ . Then, its transition kernel  $T_e$  is the identity matrix  $id_S$ . Markov Decision Process  
 089 ( $A, S, s_0, R, T$ ) is called *sparse*, if there exists a set of final states  $F \subseteq S$ , such that

$$\begin{cases} R(a, s) = 1 & T_a(s)(f) = 1 \text{ for some } f \in F \\ R(a, s) = 0 & \text{otherwise} \end{cases}$$

090 For given MDP with action space  $A$ , *trajectory* is a sequence of actions  $\alpha = a_n \dots a_1$ ; alternatively,  
 091 it's an element of free monoid on  $A$ . Given a trajectory  $\alpha$  we call its realization  $T_\alpha = T_{a_n} \cdot T_{a_{n-1}} \dots$   
 092  $T_{a_1}$ . Given the reward function  $R : A \times S \rightarrow S$ , we define  $R : A^* \rightarrow Pr(\mathbb{R})$  to be the cumulative  
 093 reward after moving along the trajectory  $\alpha \in A^*$  from the initial state  $s_0$ .

094 A *transition monoid* of MDP  $(A, S, s_0, R, T)$  is a set of trajectory realizations  $M(T) = \{T_\alpha : \alpha \in A^*\}$ , equipped with the operation of matrix multiplication. We say, that  $R$  can be *factorized through*  
 095  $M(T)$ , if there exists  $R' : M(T) \rightarrow Pr(\mathbb{R})$ , such that for all  $\alpha \in A^*$ ,  $R(\alpha) = R'(T_\alpha)$ . For brevity,  
 096 we'll write  $R$  instead of  $R'$ .

097 *Deterministic finite automaton* (DFA) is a tuple  $(Q, A, T, I, F)$ , where  $Q$  is a set of states,  $A$  -  
 098 set of actions,  $T : A \rightarrow (S \rightarrow S)$  - transition function (by  $T_a$  we'll denote a transition kernel  
 099  $T(a) : S \rightarrow S$ ),  $I$  - set of initial states,  $F$  - set of final states. Every DFA induces a directed graph  
 100 on its states. Given DFA, its *transition monoid* is a set of matrices  $\{T_\alpha : \alpha \in A^*\}$ , equipped with  
 101 the operation of matrix multiplication and identity matrix  $T_e$  as a neutral element.

### 102 3 COMPLEXITY OF MARKOV DECISION PROCESSES

103 Loosely speaking, we say that two MDPs have the same complexity if their cumulative rewards are  
 104 equal on every trajectory. We begin by proposing the extension of the notion of syntactic monoid

for general structures, such as functions, returning random variables. If one thinks that the condition  $R(a) = R(b)$  as random variables is too strict, then one could compare expectations or replace it with some approximation, for example, by  $d(R(a), R(b)) \leq \varepsilon$  after choosing certain metric on  $Pr(\mathbb{R})$  and  $\varepsilon \geq 0$ . We'll assume that the reward function of every MDP can be factorized through its transition monoid, particularly it's true for every sparse MDP.

**Definition 1.** A reward congruence  $\sim_R$  on MDP  $(A, S, s_0, R, T)$  is a congruence on its transition monoid  $\{T_\alpha : \alpha \in A^*\}$ , such that for every  $a, b \in M(T)$ ,  $a \sim_R b$  if and only if for all  $x, y \in M(T)$   $R(xay) = R(xby)$ . Then, an *irreducible monoid*  $M(R)$  of MDP is the quotient by the congruence relation  $\sim_R$ .  $R$  is well-defined on  $M(R)$ .

**Proposition 1.** A reward congruence  $\sim_R$  on  $M(T)$  for MDP  $(A, S, s_0, R, T)$  is maximal among all congruences of the type:  $\forall a, b \in M(T) a \sim b \implies R(a) = R(b)$ .

*Proof.* For a congruence  $\sim$  on  $M(T)$  and some elements  $a, b \in M(T)$ ,  $a \sim b \implies \forall x, y \in M(T) xay \sim xby$ . Hence  $R(xay) = R(xby)$ , and  $a \sim_R b$ .  $\square$

There are multiple ways to define MDP's complexity: for example, one could measure the size of state space or action space. While these definitions are reasonable, the definition we propose captures a different kind of information.

**Definition 2.** Two MDP's  $(A, S, s_0, R, T)$ ,  $(A', S', s'_0, R', T')$  are equivalent if there is an isomorphism  $\phi$  between their irreducible monoids  $M(R)$ ,  $M(R')$ , preserving reward structure, i.e.  $\forall a \in M(R), R'(\phi(a)) = R(a)$ .

The definition 2 is equivalent to another one:

**Definition 3.** Two MDP's  $(A, S, s_0, R, T)$ ,  $(A', S', s'_0, R', T')$  are *equivalent* if there is a surjective homomorphism  $\phi$  from  $A^*$   $A'^*$  or vice versa, such that reward structure is preserved:  $\forall a \in A^*, R'(\phi(a)) = R(a)$ . Hence, For deterministic MDPs with sparse binary rewards, the trajectory  $\alpha \in A^*$  solves the first MDP if and only if  $\phi(\alpha)$  solves the second MDP.

**Definition 4.** Order complexity of MDP  $(A, S, s_0, R, T)$  is the cardinality of its irreducible monoid  $M(R)$ .

**Example 1.** Suppose we want to get on the right side of the grid, which has a width of 3 and infinite length. In other words, we are given an deterministic MDP  $(A, S, s_0, R, T)$ , where:

- State space  $S = \mathbb{Z}_3 \times \mathbb{Z}$
- Initial state  $s_0 = (0, 0)$
- Action space  $A = \{(1, 0), (-1, 0), (0, 1), (0, -1)\}$
- Transition kernel  $T(i, j)(a, b) = ((a + i) \bmod 3, b + j)$
- Reward  $\begin{cases} R((i, j)(a, b)) = 1 & (i + a) \equiv 2 \pmod 3 \\ R((i, j)(a, b)) = 0 & \text{otherwise} \end{cases}$

By the definition 3, we define a reward congruence on  $M(T)$ :

$$(a, b) \sim_R (c, d) \iff \forall (x, y), (u, v), R((x, y) + (a, b) + (u, v)) = R((x, y) + (c, d) + (u, v))$$

It is true if and only if  $(x+a+u) \bmod 3 = (x+c+u) \bmod 3$ , and so  $a = c$ , since  $a, c \in \{0, 1, 2\}$ . Hence the relation  $\sim_R$  will have only 3 equivalence classes:  $\{\{(i, j) : j \in \mathbb{N}\} : i \in \mathbb{Z}_3\}$ , and the irreducible monoid will have only 3 elements. Hence, the irreducible monoid is isomorphic to  $\mathbb{Z}_3$ , and its order complexity is 3. The main conclusion is that MDP with infinite states and bigger action space might be equivalent (reduced) to MDP with three states and only one action.

We'd like to point out that even though deterministic sparse MDPs are very similar to DFAs, there exists a difference: while in RL, the episode stops if the agent reaches the final state, for DFA, the word belonging to its language might have an extension not belonging to it. Such difference can be eliminated by modifying the automaton or adding the termination action to the agent's action space. This modification looks useful and reasonable: without it, any finite MDP with a connected directed graph can be solved by the exhaustive search without any use of the environment's output.

---

162    **4 CAYLEY MAZE**  
163

164    We propose a new Open-Ended Reinforcement Learning Environment: Cayley Maze. The agent's  
165    goal is to find the path between the initial and final vertices of a directed graph by choosing the  
166    edges to move along.  
167

168    **Definition 5.** An instance of *Cayley Maze* is defined by the tuple  $(m, n, T, i, F)$ , where  
169

- 170    •
- $m \in \mathbb{N}$
- is the size of the action space , so
- $A = [m]$
- 
- 171    •
- $n \in \mathbb{N}$
- is the size of the state space, so
- $S = [n]$
- 
- 172    •
- $T : A \rightarrow \text{End}([n])$
- is the correspondence between action space and monoid generators;
- 
- 173
- $G = T(A)$
- is called the set of generators, and each generator is denoted by
- $T_a = T(a)$
- . For
- 
- 174    the function
- $T$
- extended to
- $A^*$
- :
- $T(\alpha) = T_\alpha = T(\alpha_1 \cdot \alpha_2 \dots \cdot \alpha_k) = T_{\alpha_1} \circ T_{\alpha_2} \circ \dots \circ T_{\alpha_k}$
- ,
- 
- 175    the transition monoid of
- $T$
- is the image
- $M(T) = T(A^*)$
- .
- 
- 176
- 
- 177    •
- $i \in [n]$
- is an initial state
- 
- 178    •
- $F \subseteq [n]$
- is a set of final states
- 
- 179

180    Then the instance induces a sparse deterministic MDP  $([m], [n], i, R, T)$ , where  $R$  is  
181

182    
$$\begin{cases} R(a, s) = 1 & T_a(s)(f) = 1 \text{ for some } f \in F \\ R(a, s) = 0 & \text{otherwise} \end{cases}$$
  
183

184    The opposite appears also to be true:  
185

186    **Proposition 2.** Every deterministic sparse finite MDP is an MDP of certain instance of Cayley  
187    Maze.  
188

189    *Proof.* Since  $(A, S, s_0, R, T)$  is deterministic,  $T$  can be seen as a function  $A \rightarrow (S \rightarrow S)$ . Since  
190    MDP is sparse,  $R$  is completely defined by the subset of final states  $F \subseteq S$ . Then, after enumerating  
191     $A$  and  $S$ ,  $(|A|, |S|, T, s_0, F)$  is an instance of Cayley Maze with the same MDP.  $\square$   
192

193    Cayley Mazes are parametrizable by the size of the state space  $n$ , the size of the action space  $m$ ,  
194    monoid generators  $T(A)$ , and the choice of initial and final states. From now on, all discussed  
195    Reinforcement Learning environments and their MDPs are assumed to be deterministic and sparse.  
196    While the transition between MDP and Cayley Maze formalisms is tautological, we see it valuable  
197    for several reasons.  
198

199    **4.1 CAYLEY MAZE IS NATURAL**  
200

201    Cayley Maze naturally generalizes many important problems. Such problems deserve a special  
202    name:  
203

204    **Definition 6.** An instance of the Cayley Maze is *natural* if the following holds:  
205

206    
$$(\forall \alpha, \beta \in A^* \exists i \leq n T_\alpha(i) = T_\beta(i)) \implies T_\alpha = T_\beta$$
  
207

208    Restating this property: if two trajectory realizations are equal at some state  $i$ , then they are equal at  
209    any other state. Such remarkable property can be used for evaluating agent's generalization capabili-  
210    ties and architecture's inductive biases.  
211    The problem of sorting the array of length  $n$  can be seen as the instance of natural Cayley Maze,  
212    where:  
213

- 214    • state space is the group of all
- $n$
- element permutations
- $S_n$
- ,
- 
- 215    • action space - some subset of
- $S_n$
- , generating it. For example, it could be the set of all
- 
- 216    transpositions
- $\{(i, j) : i, j \leq n\}$
- .
- 
- 217    • the transition monoid is defined by the left multiplication of the state by the action
- 
- 218

- 216     • the initial state is the unsorted number array seen as a permutation, and the final state is the  
 217       identity permutation. In this case every winning trajectory  $\alpha \in A^*$  gives the right order  $T_\alpha$   
 218

219 Rubik’s Cube is another such problem. Enumerating the squares of Rubik’s Cube allows to translate  
 220 the problem just like in the case of sorting; the only difference is that actions will have different  
 221 kinds of permutations.

#### 224 4.2 CAYLEY MAZE HAS VARIABLE COMPUTATIONAL COMPLEXITY

225 Various subfamilies of the Cayley Maze have different computational complexity: it has been shown  
 226 that the problems of sorting the array and solving the Rubik’s Cube can both be represented as  
 227 instances of the Cayley Maze. The sorting problem has polynomial time complexity, while the  
 228 problem of finding the optimal solution of Rubik’s Cube is NP-Complete (Demaine et al., 2018).  
 229

#### 230 4.3 MOST OF THE REINFORCEMENT LEARNING ENVIRONMENTS ARE NOT UNIVERSAL

232 Some of the most popular environments used to evaluate UED algorithms, like Minigrid mazes  
 233 (Boisvert et al., 2018), make heavy use of the underlying geometric structure of its state space.  
 234 We note that any Open-Ended environment, which can be solely represented by moving on the 2-  
 235 dimensional grid (i.e., for which the directed graph of its MDP can be embedded into the plane  
 236 respecting grid structure), is not universal in the sense of proposition 2: for example an MDP with  
 237 three states  $\{A, B, C\}$  and one action  $a$ :  $A \xrightarrow{a} B \xrightarrow{a} C \xrightarrow{a} A$  cannot be embedded into the plane,  
 238 since otherwise agent would have to always move in the same direction and return to the initial state.  
 239 What is less obvious is that such environments are not universal even in the sense of definition 2:

240 **Proposition 3.** There exists an MDP that is not equivalent in the sense of definition 2 to any MDP  
 241 whose directed graph is planar. Consequently, such MDP cannot be represented by moving on a  
 242 2-dimensional grid.

243 The proof of this fact is due to Book & Chandra (1976). The witnessing automaton has only 7 states  
 244 and 6 actions. The further development of this topic and the applications of topology for measuring  
 245 the complexity of finite automata can be found in Bonfante & Deloup (2018)

#### 247 4.4 MODIFICATION AND SIMPLIFICATION OF EXISTING INSTANCES

249 Given an abstract MDP or the set of game instructions, it is often unclear which modification would  
 250 make it simpler or harder. But it’s certainly possible for Cayley Mazes: painting all Rubik’s cube  
 251 faces into black makes it much easier to solve. In other words, given the MDP  $M$  whose transition  
 252 monoid acts on the set  $S$ , and the coloring of  $S$  - a surjective function  $h : S \rightarrow K$ , it’s not hard to  
 253 build *quotient* of  $M$  by  $h$ , whose construction is similar to its group-theoretic analog.

254 Many constructions on groups, such as products, allow to efficiently combine existing MDPs to  
 255 produce new ones.

### 257 5 APPLICATIONS AND IMPLEMENTATION DETAILS

259 We implement Cayley Maze as a parametrizable environment in JAX with a Gym interface. It allows  
 260 the sampling of any MDP with predefined action space and state spaces or the modification of the  
 261 existing transition monoid, initial and target states.

262 Cayley Maze may be used in the evaluation procedure of every learning problem, which has a  
 263 sequential structure and where generalization or emergent complexity of observations are crucial.  
 264 Some of the proposed scenarios are:

- 265     1. Since universality is essential for UED, Cayley Maze may be used for the evaluation of  
 266       UED algorithms  
 267  
 268     2. It is possible to create environment samplers whose instances have a common structure. For  
 269       example, it might be MDPs whose transition monoids are simple groups or environments  
 270       that represent  $n \times n \times n$  Rubik’s Cube.

- 270        3. It is possible to evaluate not only on the subfamilies of environments but also on various  
 271        combinations of these environments. For example, we can check whether the agent, which  
 272        performed well on some instances, would perform well on its product.  
 273  
 274        4. Another way to test generalization capabilities is to use the local property of natural Cayley  
 275        Mazes, as in 4.1. For example, evaluate the agent's performance on the initial states that  
 276        were unreachable during training.  
 277  
 278        5. Since the process of creating a new MDP can be seen as an MDP, it can be expressed as an  
 279        instance of Cayley Maze. Hence, the UED scenarios where the teacher learns to build an  
 environment without student become possible.

280        Another interesting feature of the implementation is that while constructing the most general re-  
 281        inforcement learning environment, one might expect the explosion of the teacher's action space.  
 282        The implementation allows customization of the desired trade-off between the size of action space,  
 283        representation dimension, and the power of the edit per step.

## 285        REFERENCES

- 287        M. Beukman, S. Coward, et al. Refining minimax regret for unsupervised environment design, 2024.  
 288        URL <https://arxiv.org/abs/2402.12284>.
- 289        C. Boisvert, L. Willems, et al. Minimalistic gridworld environment for ope- nai gym, 2018. URL  
 290        <https://github.com/maximecb/gym-minigrid>.
- 292        G. Bonfante and F. Deloup. The genus of regular languages. *Mathematical Structures in Computer  
 293        Science*, 28(1):14–44, 2018.
- 294        R. Book and A. Chandra. Inherently nonplanar automata. *A.K. Acta Informatica*, 6:89–94, 1976.
- 296        J. Bruce, M. Dennis, et al. Genie: Generative interactive environments, 2024. URL <https://arxiv.org/abs/2402.15391>.
- 299        E. Demaine, S. Eisenstat, et al. Solving the rubik's cube optimally is np-complete. In *35th Sym-  
 300        posium on Theoretical Aspects of Computer Science (STACS 2018). Leibniz International Proceed-  
 301        ings in Informatics (LIPIcs)*, 96:24:1–24:13, 2018.
- 302        E. Hughes, M. Dennis, et al. Open-endedness is essential for artificial superhuman intelligence,  
 303        2024. URL <https://arxiv.org/abs/2406.04268>.
- 304        M. Matthews, M. Beukman, et al. Craftax: A lightning-fast benchmark for open-ended reinforce-  
 305        ment learning, 2024. URL <https://arxiv.org/abs/2402.16801>.
- 307        J. Parker-Holder, M. Jiang, et al. Evolving curricula with regret-based environment design, 2023.  
 308        URL <https://arxiv.org/abs/2203.01302>.
- 309        J. Pin. *Handbook of automata theory. Volume I. Theoretical foundations*, volume 1. Berlin: Euro-  
 310        pean Mathematical Society (EMS), 2021.