

503 **Appendix A Videos**

- 504 - Video 1 ([https://drive.google.com/open?id=1aVL7fmNp\\_rMpU2E00PXq6BGtDNvmvCTu](https://drive.google.com/open?id=1aVL7fmNp_rMpU2E00PXq6BGtDNvmvCTu))  
 505 - Video 2 ([https://drive.google.com/open?id=1s16Mz66ETV\\_dpLXi0yyaJLJHZpFN01uv](https://drive.google.com/open?id=1s16Mz66ETV_dpLXi0yyaJLJHZpFN01uv))

506 **Appendix B Extra Figures**

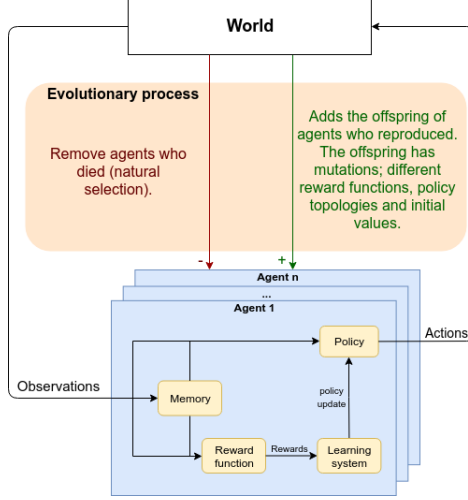


Figure 3: A diagram on how evolution and learning affect the development of the brain.

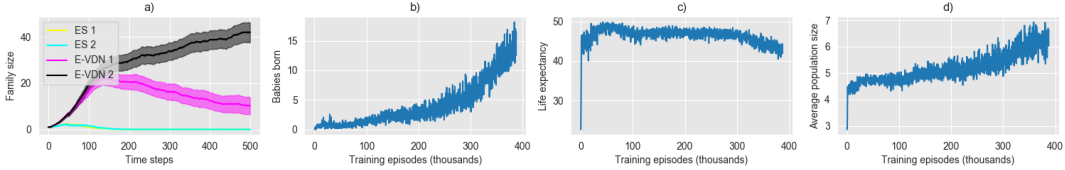


Figure 4: a) (CMA-ES vs E-VDN) Average and 95% confidence interval of two CMA-ES and two E-VDN family sizes computed over 90 episodes. b, c and d) the macro-statistics obtained in the non-binary environment. To speed up the training we used the smaller NN and the denser evolutionary reward described in the Appendix G.3.

507 **Appendix C Optimally Aligned Reward: Formal Description**

508 The notion of an optimal reward function for a given fitness function was introduced by Singh [21, 33],  
 509 here we adapt his original formulation. From the perspective of agent  $i$  the environment is defined by  
 510 the state transition distribution;  $E^i := p(s_{t+1}^i | s_t^i, a_t^i, \pi^{-i})$ . Where  $\pi^{-i}$  is the concatenation of the  
 511 policies of all agents except agent  $i$ ;  $\pi^{-i} := \{\pi^j\}_{j \neq i}$ . Formally we say:  $h_\pi^i \sim \langle \mathcal{L}(\mathcal{R}^i), E^i \rangle$ , where  
 512  $h_\pi^i$  is the sampled history of the adaptations of policy  $\pi^i$  which resulted from agent  $i$  learning  $\mathcal{L}(\cdot)$  to  
 513 maximise its reward function  $\mathcal{R}^i$  by interacting  $\langle \cdot \rangle$  with the environment. If agent  $i$  is the only agent  
 514 learning then  $\pi^{-i}$  is static and so is the environment  $E^i$ . In this case, the optimally aligned reward  
 515 function is given by:

$$\mathcal{R}^* = \arg \max_{\mathcal{R}^i} \mathbb{E}_{h_\pi^i \sim \langle \mathcal{L}(\mathcal{R}^i), E^i \rangle} \mathcal{F}(h_\pi^i, E^i), \quad (10)$$

516 where  $\mathcal{F}$  is the fitness function. In the general case, all agents are learning, and therefore, the  
 517 environment is non-static, the fitness for  $h_\pi^i$  is changing and so is the optimally aligned reward  $\mathcal{R}^*$ .

## 518 Appendix D Value-Decomposition Networks

519 Our work builds on VDN [38], which was designed to address the binary cooperative MARL setting.  
 520 In this setting, the agents are grouped into teams and all the agents within a team receive the  
 521 same reward. VDN’s main assumption is that the joint action-value function of the whole team of  
 522 cooperative agents can be additively decomposed into the action-value functions across the members  
 523 of the team.

$$Q^{\mathcal{T}}((h_t^1, h_t^2, \dots, h_t^{|\mathcal{T}|}), (a_t^1, a_t^2, \dots, a_t^{|\mathcal{T}|})) \approx \sum_{i \in \mathcal{T}} \tilde{Q}^i(h_t^i, a_t^i), \quad (11)$$

524 where  $\mathcal{T}$  is the set of agents belonging to the team, and  $\tilde{Q}^i(h_t^i, a_t^i)$  is the value function of agent  $i$   
 525 which depends solely on its partial observation of the environment and its action at time  $t$ .  $\tilde{Q}^i$  are  
 526 trained by back-propagating gradients from the Q-learning rule through the summation.

$$g_i = \nabla \theta_i (y_t^{\mathcal{T}} - \sum_{i \in \mathcal{T}} \tilde{Q}(h_t^i, a_t^i | \theta_i))^2, \quad y_t^{\mathcal{T}} = r_t^{\mathcal{T}} + \gamma \sum_{i \in \mathcal{T}} \max_{a_{t+1}^i} \tilde{Q}(h_{t+1}^i, a_{t+1}^i | \theta_i), \quad (12)$$

527 where  $\theta_i$  are the parameters of  $\tilde{Q}^i$ ,  $g_i$  is its gradient and  $r_t^{\mathcal{T}}$  is the reward for the team  $\mathcal{T}$  at the time  
 528 instant  $t$ . Note that even though the training process is centralised, the learned agents can be deployed  
 529 independently, since each agent acting greedily with respect to its own  $\tilde{Q}^i$  will also maximise its  
 530 team value function  $\arg \max_{a_t^i} Q_t^{\mathcal{T}}(\dots) \approx \arg \max_{a_t^i} \tilde{Q}^i(h_t^i, a_t^i)$ .

## 531 Appendix E Environment

532 In this section, we go through the game loop of the environments summarised in the main article.  
 533 Both the binary and non-binary environment have the same game loop, their only difference is in the  
 534 way agents reproduce and in the length of the genome agents carry (the genome has a single gene in  
 535 the binary environment and 32 genes in the non-binary one). The states of the tiles and agents are  
 536 described in table 1.

	Tile state	Agent state	
Type	Boolean (food source/dirt)	Position (x,y)	Integer, Integer
Occupied	Boolean	Health	Integer
Food available	Float	Age	Integer
		Food stored	Float
		Genome	Integer Vector

Table 1: The state of the tiles and agents.

537 We now introduce the various components of the game loop:

538 **Initialisation** The simulation starts with five agents, each one with a unique genome. All agents  
 539 start with age 0 and  $e$  units of food (the endowment). The environments are never-ending. Table 2  
 540 describes the configuration used in the paper.

Endowment ( $e$ )	Initial health	Start of fertility age	End of fertility age	Longevity	World size	Food growth rate ( $f_r$ )	Maximum food capacity ( $c_f$ )
10	2	5	40	50	50x50	0.15	3

Table 2: Configuration of the environment used in the paper.

541 **Food production** Each tile on the grid world can either be a food source or dirt. Food sources  
 542 generate  $f_r$  units of food per iteration until reaching their maximum food capacity ( $c_f$ ).

543 **Foraging** At each iteration, an agent can move one step to North, East, South, West or choose to  
544 remain still. When an agent moves to a tile with food it collects all the available food in it. The map  
545 boundaries are connected (e.g. an agent that moves over the top goes to the bottom part of the map).  
546 Invalid actions, like moving to an already occupied tile, are ignored.

547 **Attacking** At each iteration, an agent can also decide to attack a random adjacent agent: this is an  
548 agent within one step to N, E, S or W. Each attack takes 1 unit of health from the victim's. If the  
549 victim's health reaches zero, it dies, and the attacker will "eat it" and receive 50% of its food reserves.

550 **Asexual Reproduction** An agent is considered fertile if it has accumulated more than twice the  
551 amount of food it received at birth (i.e. twice its endowment  $e$ ) and its age is within a given fertile  
552 age. The fertile agent will give birth once they have an empty tile nearby, when that happens the  
553 parent transfers  $e$  units of food to its newborn child. The newborn child will have the same genome  
554 has its parent.

555 **Sexual Reproduction** An agent is considered fertile if it has accumulated more than the amount of  
556 food it received at birth and its age is within a given fertile age. The fertile agent will give birth once  
557 it is adjacent to another fertile agent and one of them has an empty tile nearby, when that happens  
558 each parent transfers  $\frac{e}{2}$  units of food to its newborn child. A random half of the newborn's genes  
559 come from the first parent, and the second half comes from the second parent.

560 **Game loop** At every iteration, we randomise the order at which the agents execute their actions.  
561 Only after all the agents are in their new positions, the attacks are executed (with the same order as  
562 the movement actions). The complete game loop is summarized in the next paragraph.

563 At each iteration, each agent does the following:

- 564 • Execute a movement action: Stay still or move one step North, East, South or West.
- 565 • Harvest: Collect all the food contained in its tile.
- 566 • Reproduce: Give birth to a child using asexual or sexual reproduction (see their respective  
567 sections).
- 568 • Eat: Consume a unit of food.
- 569 • Age: Get one year older.
- 570 • Die: If an agent's food reserves become empty or it becomes older than its longevity  
571 threshold, then it dies.
- 572 • Execute an attack action: After every agent has moved, harvested, reproduced, eaten and  
573 aged the attacks are executed. Agents that reach zero health get eaten at this stage.

574 Additionally, at each iteration, each food source generates  $f_r$  units of food until reaching the given  
575 maximum capacity ( $c_f$ ).

## 576 **Appendix F Evolution Strategies**

577 In the binary environment, we compare the E-VDN algorithm with a popular ES algorithm. ES  
578 algorithms optimise an agent's policy by sampling policy weights from a multivariate Gaussian  
579 distribution, evaluating those weights on the environment, giving them a fitness score and updating  
580 the Gaussian distribution parameters so that the next samples are more likely to achieve a higher  
581 fitness. There are a few different methods on how to update the distribution parameters, we chose to  
582 use CMA-ES [14] because it has been successful in optimising NN for a wide range of sequential  
583 decision problems [16, 17, 18]. However, note that CMA-ES was not designed for multi-agent  
584 settings where the fitness function landscape changes as the other agents learn. Nevertheless, we  
585 used five independent multivariate Gaussian distributions each one associated with a unique gene  
586 and each one being updated by the CMA-ES algorithm. In the beginning, when the agents can not  
587 survive for long, the fitness function is given by the total sum of family members along time, once  
588 the agents learn how to survive and reproduce we change the fitness function to be the number of  
589 family members at the end of an episode with 500 steps. Since the CMA-ES algorithm computation

590 time grows quadratic with the number of parameters,  $O(N^2)$ , we had to use the smaller NN for this  
 591 comparison. The algorithm was implemented using an available *python* library [15].

592 In the non-binary environment, if the initial five agents have different policies it creates the problem  
 593 of deciding which policy should a child inherit. At the same time, the initial five agents can't all share  
 594 the same policy because it then becomes impossible to define the fitness function for each policy. For  
 595 these reasons, we didn't implement CMA-ES on the non-binary environment.

## 596 Appendix G Algorithm details

### 597 G.1 Effective time horizon

We want to find the number of iterations ( $h_e$ ) that guarantee an error between the estimate of the final  
 reward and the actual final reward to be less or equal than a given  $\epsilon$ ,  $|r_{T^i-1}^i - \hat{r}_{T^i-1}^i| \leq \epsilon$ .  
 Remember that the final reward is given by:

$$r_{T^i-1}^i = \sum_{t=T^i}^{\infty} \gamma^{t-T^i} \sum_{j \in \mathcal{A}_t} k(\mathbf{g}^i, \mathbf{g}^j) = \sum_{t'=0}^{\infty} \gamma^{t'} k_{t'}^i$$

598 Where  $t' = t - T^i$  and  $k_{t'}^i = \sum_{j \in \mathcal{A}_{t'}} k(\mathbf{g}^i, \mathbf{g}^j)$ . The estimate of the final reward is computed with  
 599 the following finite sum  $\hat{r}_t^i = \sum_{t'=0}^{h_e-1} \gamma^{t'} k_{t'}^i$ .  
 600

601 Note that  $k_t^i$  is always positive so the error  $r_{T^i-1}^i - \hat{r}_{T^i-1}^i$  is always positive as well. To find the  $h_e$   
 602 that guarantees an error smaller or equal to epsilon we define  $r_b$  as the upper bound of  $k_t^i$  and ensure  
 603 that the worst possible error is smaller or equal to epsilon:

$$\sum_{t'=0}^{\infty} \gamma^{t'} r_b - \sum_{t'=0}^{h_e-1} \gamma^{t'} r_b \leq \epsilon \quad (13)$$

$$\frac{r_b}{1-\gamma} - r_b \frac{1-\gamma^{h_e}}{1-\gamma} \leq \epsilon \quad (14)$$

$$\frac{r_b \gamma^{h_e}}{1-\gamma} \leq \epsilon \quad (15)$$

$$h_e \log \gamma \leq \log \frac{\epsilon(1-\gamma)}{r_b} \quad (16)$$

$$h_e \leq \frac{\log \frac{\epsilon(1-\gamma)}{r_b}}{\log \gamma} \quad (17)$$

604 We go from (1) to (2) by using the known convergences of geometric series:  $\sum_{k=0}^{\infty} ar^k = \frac{a}{1-r}$   
 605 and  $\sum_{k=0}^{n-1} ar^k = a \frac{1-r^n}{1-r}$  for  $r < 1$ . Since  $h_e$  needs to be a positive integer we take the ceil  
 606  $h_e = \left\lceil \frac{\log \frac{\epsilon(1-\gamma)}{r_b}}{\log \gamma} \right\rceil$  and note that this equation is only valid when  $\frac{\epsilon(1-\gamma)}{r_b} < 1$ . For example, an  
 607 environment that has the capacity to feed at most 100 agents has an  $r_b = 100$  (which is the best  
 608 possible reward, i.e. the kinship between every agent is 1). If we use  $\epsilon = 0.1$  and  $\gamma = 0.9$  then  
 609  $h_e = 88$ .

### 610 G.2 Experience buffer

611 When using Q-learning methods with DQN, as we are, it's common practice to use a replay buffer.  
 612 The replay buffer stores the experiences  $(s_t, a_t, r_t, s_{t+1})$  for multiple time steps  $t$ . When training, the  
 613 algorithm randomly samples experiences from the replay buffer. This breaks the auto-correlation  
 614 between the consecutive examples and makes the algorithm more stable and sample efficient. How-  
 615 ever, for non-stationary environments, past experiences might be outdated. For this reason, we don't  
 616 use a replay buffer. Instead, we break the auto-correlations by collecting experiences from many  
 617 independent environments being sampled in parallel. After a batch of experiences is used we discard

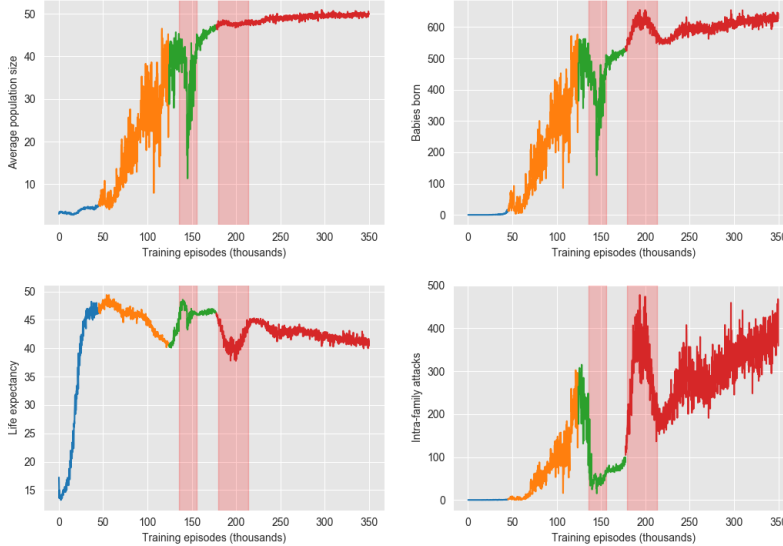


Figure 5: Macro-statistics when evolving bacteria using the standard evolutionary reward. I: learning to survive (blue line), II: learning to reproduce (orange), III: learning to detect kin (green), IV: learning to self-sacrifice (red). The red bands correspond to the First and Second Family Wars.

618 them. In our experiments, we simulated 400 environments in parallel and collected one experience  
 619 step from each agent at each environment to form a training batch.

### 620 G.3 Denser reward function

621 In some situations, we used a denser version of the evolutionary reward to speed up the training  
 622 process. We call it the *sugary* reward,  $r_t^i = \sum_{j \in \mathcal{A}_t} k(\mathbf{g}^i, \mathbf{g}^j) f_t^j$  where  $f_t^j$  is the food collected by  
 623 agent  $j$  at the time instant  $t$ . In these simple environments, the *sugary* and the evolutionary reward are  
 624 almost equivalent since a family with more members will be able to collect more food and vice-versa.  
 625 However, the *sugary* reward contains more immediate information whilst the evolutionary reward has  
 626 a lag between good (bad) actions and high (low) rewards; a family that is not doing a good job at  
 627 collecting food will take a while to see some of its members die from starvation. Nonetheless, the  
 628 evolutionary reward is more correct since it describes exactly what we want to maximise. Note that  
 629 this reward was not used to produce the results when comparing E-VDN with CMA-ES.

630 When using the standard evolutionary reward to evolve the larger NNs, the same four eras, that were  
 631 observed with the *sugary* reward, emerge. However, their progression is not as linear. In this case, the  
 632 families take longer to learn and sometimes one family evolves much faster than the others. When this  
 633 happens, the families left behind eventually catch up with the most developed ones. The behaviour of  
 634 the emerging families successfully interferes with the developed ones creating a temporary disruption  
 635 in the environment which disrupts its macro-statistics. Two disruptions were observed in one of our  
 636 simulations and we named them the First and the Second Family Wars (fig. 5).

## 637 Appendix H Results details

### 638 H.1 Cannibalism and suicide as a tool for gene survival

639 In the evolutionary history of the binary environment we saw the rise of cannibalism in the fourth era.  
 640 Figure 6.a shows how the average age of cannibals and their victims grows apart in this era. After  
 641 observing this behaviour, we wanted to know how important cannibalism was for gene survival. To  
 642 answer this question, we measured the family size of a certain family when its members were not  
 643 allowed to attack each other and compared it with the normal situation where intra-family attacks  
 644 were allowed (see Figure 6.b). Figure 6.b clearly shows that, in this environment, cannibalism is  
 645 essential for long-term gene survival. We also ran this exact experiment before the fourth era and

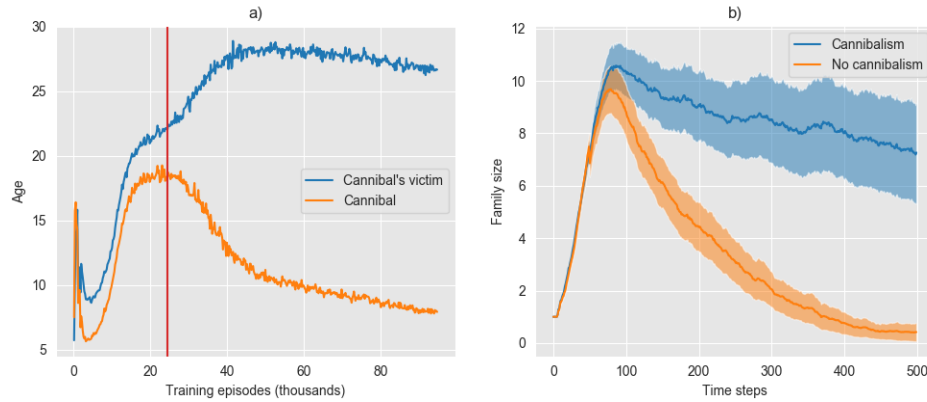


Figure 6: a) The average age of intra-family cannibals and cannibals' victims. The vertical red line marks the start of era IV. b) the size of family 1 averaged along 90 test episodes. To compute the orange line we simply blocked all the attacks between members of the family 1. The shaded bands represent the 95% confidence interval.

646 achieved the opposite results, suggesting that before this era the agents didn't yet know how to use  
 647 cannibalism to their gene's advantage (results not shown).

## 648 H.2 Genetic drift

649 The environment starts with a unique set of alleles<sup>1</sup>, and there is no mechanism to add diversity to  
 650 this initial set (no mutations) but there is a mechanism to remove diversity: death. When a death  
 651 occurs, there is a chance that the last carrier of a particular allele is lost and if there are no mutations,  
 652 this is a permanent loss in diversity. This evolutionary mechanism that changes allele frequencies  
 653 using chance events is called genetic drift. In the binary environment, this means that if we simulate  
 654 the environment for long enough all the agents will end up sharing the same allele. In the non-binary  
 655 environment, all the agents will have the same genome, however, this genome will likely be composed  
 656 by alleles coming from all the five founders (see Video 2).

657 We found that in the binary environment kin detection speeds up this decrease in allele diversity.  
 658 This was expected since agents cooperate with kin and compete with non-kin. Therefore, as a family  
 659 gets bigger, its members become more likely to encounter cooperative family members rather than  
 660 competitive unrelated agents. This improves the survival and reproduction success of that family,  
 661 making it even bigger. Figure 7 shows the decrease in diversity with and without kin selection (we  
 662 removed kin selection by zeroing out the kinship feature in the agents' observations). From the figure,  
 663 it is evident that kin selection speeds up the decrease in diversity. However, note that before the  
 664 100<sup>th</sup> iteration, kin detection leads to a slightly higher diversity. This happens because kin detection  
 665 reduces intra-family violence, leading to fewer deaths and consequently to a slower genetic drift.  
 666 The environment usually reaches its maximum capacity around the 100<sup>th</sup> iteration, at this time the  
 667 inter-family competition is at its highest and the positive feedback loop created by kin selection starts  
 668 having a larger importance.

## 669 References

- 670 [1] David Ackley and Michael Littman. Interactions between learning and evolution. *Artificial life*  
 671 *II*, 10:487–509, 1991.
- 672 [2] BURT Austin, Robert Trivers, and Austin Burt. *Genes in conflict: the biology of selfish genetic*  
 673 *elements*. Harvard University Press, 2009.

<sup>1</sup>A gene can have many alleles — a variant form of the same gene. For example, a gene that codes the eye colour may have two different alleles one for brown eyes and other for blue eyes. In everyday language, the term gene is often used when referring to an allele, when one says that two siblings carry the same eye colour gene they actually mean that they carry the same allele. We will use this language as well, but it should be clear from the context to which one we are referring.

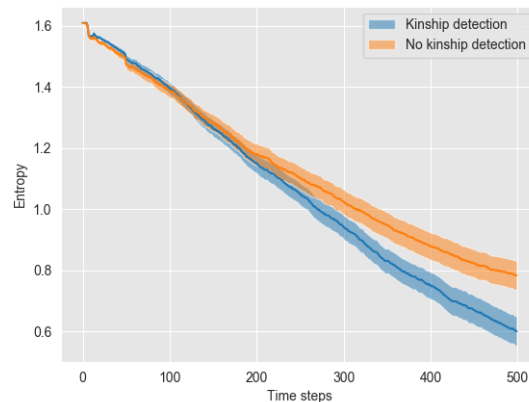


Figure 7: Entropy (diversity) in the allele frequency along the first 500 time steps in the situation where agents can detect kinship and when they can't. The bands show the 95% confidence interval.

- 674 [3] Trapit Bansal, Jakub Pachocki, Szymon Sidor, Ilya Sutskever, and Igor Mordatch. Emergent  
675 complexity via multi-agent competition. *arXiv preprint arXiv:1710.03748*, 2017.
- 676 [4] Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemysław Dębiak, Christy  
677 Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, et al. Dota 2 with large  
678 scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680*, 2019.
- 679 [5] Wendelin Böhmer, Vitaly Kurin, and Shimon Whiteson. Deep coordination graphs. *arXiv  
680 preprint arXiv:1910.00091*, 2019.
- 681 [6] Richard Dawkins. *The selfish gene*. 1976.
- 682 [7] Richard Dawkins. Twelve misunderstandings of kin selection. *Zeitschrift für Tierpsychologie*,  
683 51(2):184–200, 1979.
- 684 [8] Stefan Elfving, Eiji Uchibe, Kenji Doya, and Henrik I Christensen. Co-evolution of shaping  
685 rewards and meta-parameters in reinforcement learning. *Adaptive Behavior*, 16(6):400–412,  
686 2008.
- 687 [9] Jakob N Foerster, Gregory Farquhar, Triantafyllos Afouras, Nantas Nardelli, and Shimon  
688 Whiteson. Counterfactual multi-agent policy gradients. In *Thirty-Second AAAI Conference on  
689 Artificial Intelligence*, 2018.
- 690 [10] Andy Gardner and Francisco Úbeda. The meaning of intragenomic conflict. *Nature ecology &  
691 evolution*, 1(12):1807–1815, 2017.
- 692 [11] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil  
693 Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural  
694 information processing systems*, pages 2672–2680, 2014.
- 695 [12] Carlos Guestrin, Daphne Koller, and Ronald Parr. Multiagent planning with factored mdps. In  
696 *Advances in neural information processing systems*, pages 1523–1530, 2002.
- 697 [13] Eric A Hansen, Daniel S Bernstein, and Shlomo Zilberstein. Dynamic programming for partially  
698 observable stochastic games. In *AAAI*, volume 4, pages 709–715, 2004.
- 699 [14] Nikolaus Hansen, Sibylle D Müller, and Petros Koumoutsakos. Reducing the time complexity of  
700 the derandomized evolution strategy with covariance matrix adaptation (cma-es). *Evolutionary  
701 computation*, 11(1):1–18, 2003.
- 702 [15] Nikolaus Hansen, Youhei Akimoto, and Petr Baudis. CMA-ES/pycma on Github. Zenodo,  
703 DOI:10.5281/zenodo.2559634, February 2019. URL [https://doi.org/10.5281/zenodo.  
704 2559634](https://doi.org/10.5281/zenodo.2559634).

- 705 [16] Verena Heidrich-Meisner and Christian Igel. Evolution strategies for direct policy search. In  
706 *International Conference on Parallel Problem Solving from Nature*, pages 428–437. Springer,  
707 2008.
- 708 [17] Verena Heidrich-Meisner and Christian Igel. Variable metric reinforcement learning methods  
709 applied to the noisy mountain car problem. In *European Workshop on Reinforcement Learning*,  
710 pages 136–150. Springer, 2008.
- 711 [18] Verena Heidrich-Meisner and Christian Igel. Hoeffding and bernstein races for selecting policies  
712 in evolutionary direct policy search. In *Proceedings of the 26th Annual International Conference*  
713 *on Machine Learning*, pages 401–408. ACM, 2009.
- 714 [19] Jelle R Kok, Eter Jan Hoen, Bram Bakker, and Nikos Vlassis. Utile coordination: Learning  
715 interdependencies among cooperative agents. In *EEE Symp. on Computational Intelligence and*  
716 *Games, Colchester, Essex*, pages 29–36, 2005.
- 717 [20] Karol Kurach, Anton Raichuk, Piotr Stańczyk, Michał Zając, Olivier Bachem, Lasse Espeholt,  
718 Carlos Riquelme, Damien Vincent, Marcin Michalski, Olivier Bousquet, et al. Google research  
719 football: A novel reinforcement learning environment. *arXiv preprint arXiv:1907.11180*, 2019.
- 720 [21] Richard L Lewis, Satinder Singh, and Andrew G Barto. Where do rewards come from? In  
721 *Proceedings of the International Symposium on AI-Inspired Biology*, pages 2601–2606, 2010.
- 722 [22] Ryan Lowe, Yi Wu, Aviv Tamar, Jean Harb, OpenAI Pieter Abbeel, and Igor Mordatch. Multi-  
723 agent actor-critic for mixed cooperative-competitive environments. In *Advances in Neural*  
724 *Information Processing Systems*, pages 6379–6390, 2017.
- 725 [23] Simon M Lucas and Thomas P Runarsson. Temporal difference learning versus co-evolution for  
726 acquiring othello position evaluation. In *2006 IEEE Symposium on Computational Intelligence*  
727 *and Games*, pages 52–59. IEEE, 2006.
- 728 [24] Simon M Lucas and Julian Togelius. Point-to-point car racing: an initial study of evolution  
729 versus temporal difference learning. In *2007 IEEE symposium on computational intelligence*  
730 *and games*, pages 260–267. IEEE, 2007.
- 731 [25] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G  
732 Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al.  
733 Human-level control through deep reinforcement learning. *Nature*, 518(7540):529, 2015.
- 734 [26] Frans A. Oliehoek and Christopher Amato. *A Concise Introduction to Decentralized POMDPs*.  
735 SpringerBriefs in Intelligent Systems. Springer, May 2016. doi: 10.1007/978-3-319-28929-8.  
736 URL <http://www.fransoliehoek.net/docs/OliehoekAmato16book.pdf>.
- 737 [27] Frans A Oliehoek, Shimon Whiteson, Matthijs TJ Spaan, et al. Approximate solutions for  
738 factored dec-pomdps with many agents. In *AAMAS*, pages 563–570, 2013.
- 739 [28] Tabish Rashid, Mikayel Samvelyan, Christian Schroeder De Witt, Gregory Farquhar, Jakob  
740 Foerster, and Shimon Whiteson. Qmix: monotonic value function factorisation for deep  
741 multi-agent reinforcement learning. *arXiv preprint arXiv:1803.11485*, 2018.
- 742 [29] Craig W Reynolds. Competition, coevolution and the game of tag. In *Proceedings of the Fourth*  
743 *International Workshop on the Synthesis and Simulation of Living Systems*, pages 59–69, 1994.
- 744 [30] Thomas Philip Runarsson and Simon M Lucas. Coevolution versus self-play temporal difference  
745 learning for acquiring position evaluation in small-board go. *IEEE Transactions on Evolutionary*  
746 *Computation*, 9(6):628–640, 2005.
- 747 [31] David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur  
748 Guez, Marc Lanctot, Laurent Sifre, Dharshan Kumaran, Thore Graepel, et al. A general  
749 reinforcement learning algorithm that masters chess, shogi, and go through self-play. *Science*,  
750 362(6419):1140–1144, 2018.
- 751 [32] Karl Sims. Evolving 3d morphology and behavior by competition. *Artificial life*, 1(4):353–372,  
752 1994.



- 753 [33] Satinder Singh, Richard L Lewis, Andrew G Barto, and Jonathan Sorg. Intrinsically motivated  
754 reinforcement learning: An evolutionary perspective. *IEEE Transactions on Autonomous Mental*  
755 *Development*, 2(2):70–82, 2010.
- 756 [34] J Maynard Smith. Group selection and kin selection. *Nature*, 201(4924):1145, 1964.
- 757 [35] J Maynard Smith and George R Price. The logic of animal conflict. *Nature*, 246(5427):15,  
758 1973.
- 759 [36] Kenneth O Stanley and Risto Miikkulainen. Evolving neural networks through augmenting  
760 topologies. *Evolutionary computation*, 10(2):99–127, 2002.
- 761 [37] Joseph Suarez, Yilun Du, Phillip Isola, and Igor Mordatch. Neural mmo: A massively  
762 multiagent game environment for training and evaluating intelligent agents. *arXiv preprint*  
763 *arXiv:1903.00784*, 2019.
- 764 [38] Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinicius Zambaldi,  
765 Max Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z Leibo, Karl Tuyls, et al. Value-  
766 decomposition networks for cooperative multi-agent learning. *arXiv preprint arXiv:1706.05296*,  
767 2017.
- 768 [39] Richard S Sutton, Andrew G Barto, et al. *Introduction to reinforcement learning*, volume 2.  
769 MIT press Cambridge, 1998.
- 770 [40] Elise Van der Pol and Frans A Oliehoek. Coordinated deep reinforcement learners for traffic  
771 light control. *Proceedings of Learning, Inference and Control of Multi-Agent Systems (at NIPS*  
772 *2016)*, 2016.
- 773 [41] Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki, Michaël Mathieu, Andrew Dudzik, Jun-  
774 young Chung, David H Choi, Richard Powell, Timo Ewalds, Petko Georgiev, et al. Grandmaster  
775 level in starcraft ii using multi-agent reinforcement learning. *Nature*, 575(7782):350–354, 2019.
- 776 [42] Shimon Whiteson and Peter Stone. Evolutionary function approximation for reinforcement  
777 learning. *Journal of Machine Learning Research*, 7(May):877–917, 2006.
- 778 [43] George C Williams. *Adaptation and natural selection: a critique of some current evolutionary*  
779 *thought*, volume 833082108. Princeton science library OCLC, 1966.