

SIMPLIFIED ACTION DECODER FOR DEEP MULTI-AGENT REINFORCEMENT LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

In recent years we have seen fast progress on a number of benchmark problems in AI, with modern methods achieving near or super human performance in Go, Poker and Dota. One common aspect of all of these challenges is that they are by design *adversarial* or, technically speaking, zero-sum. In contrast to these settings, success in the real world commonly requires humans to collaborate and communicate with others, in settings that are, at least partially, cooperative. In the last year, the card game *Hanabi* has been established as a new benchmark environment for AI to fill this gap. In particular, Hanabi is interesting to humans since it is entirely focused on *theory of mind*, *i.e.*, the ability to effectively reason over the intentions, beliefs and point of view of other agents when observing their actions. Learning to be informative when observed by others is an interesting challenge for Reinforcement Learning (RL): Fundamentally, RL requires agents to explore in order to discover good policies. However, when done naively, this randomness will inherently make their actions less informative to others during training. We present a new deep multi-agent RL method, the *Simplified Action Decoder* (SAD), which resolves this contradiction exploiting the centralized training phase. During training SAD allows other agents to not only observe the (*exploratory*) action chosen, but agents instead also observe the *greedy* action of their team mates. By combining this simple intuition with an auxiliary task for state prediction and best practices for multi-agent learning, SAD establishes a new state of the art for 2-5 players on the self-play part of the Hanabi challenge.

1 INTRODUCTION

Humans are highly social animals and spend vast amounts of time coordinate, collaborating and communicate with others. In contrast to these, at least partially, cooperative multi-agent settings, most progress on AI in games has been in zero-sum games where agents compete against each other, typically rendering communication futile. This includes examples such as Go (Silver et al., 2016; 2017; 2018), poker (Brown & Sandholm, 2017; Moravčík et al., 2017; Brown & Sandholm, 2019) and chess (Campbell et al., 2002).

This narrow focus is unfortunate, since communication and coordination require unique abilities. In order to enable smooth and efficient social interactions of groups of people, it is commonly required to reason over the intents, points of views and beliefs of other agents from observing their actions. For example, a driver can reasonably infer that if a truck in front of them is slowing down when approaching an intersection, then there is likely an obstacle ahead. Furthermore, humans are both able to interpret the actions of others and can act in a way that is informative when their actions are being observed by others, capabilities that are commonly called *theory of Mind* (ToM), (Baker et al., 2017). Importantly, in order to carry out this kind of reasoning, an agent needs to consider why a given action is taken and what this decision indicates about the state of the world. Simply observing other agents are doing is not sufficient.

While these kind of abilities are particularly relevant in partially observable fully cooperative multi-agent settings, ToM reasoning clearly matters in a variety of real world scenarios. For example, autonomous cars will likely need to understand the point of view, intents and beliefs of other traffic participants in order to deal with highly interactive settings such as 4-way crossing or dense traffic in cities.

Hanabi is a fully cooperative partially-observable card game that has recently been proposed as a new benchmark challenge problem for AI research (Bard et al., 2019) to fill the gap around ToM. In Hanabi, players need to find conventions that allow them to effectively exchange information from their local observations through their actions, taking advantage of the fact that actions are observed by all team mates.

Most prior state-of-the-art agents for Hanabi were developed using handcrafted algorithms, which beat off-the-shelf deep multi-agent RL methods by a large margin. This makes intuitive sense: Beyond the “standard” multi-agent challenges of credit assignment, non-stationarity and joint exploration, learning an informative policy presents an additional fundamentally new conflict. On the one hand, an RL agent needs to explore in order to discover good policies through trial and error. On the other hand, when carried out naively, this exploration will add noise to the policy of the agent during the training process, making their actions strictly less informative to their team mates.

One possible solution to this is to explore in the space of deterministic partial policies, rather than actions, and sample these policies from a distribution that conditions on a *common knowledge* Bayesian belief. This is successfully carried out in the Bayesian Action Decoder (BAD) (Foerster et al., 2019), the only previous Deep RL method to accomplish a state of the art in Hanabi. While this is a notable accomplishment, it comes at the cost of simplicity and generality. For a start, BAD requires an explicit common knowledge Bayesian belief to be tracked, which not only adds computational burden due to the required sampling steps, but also uses expert knowledge regarding the game dynamics. Furthermore, BAD, as presented, is trained using actor-critic methods which are sample inefficient and suffer from local optima. In order to get around this, BAD uses population based training, further increasing the number of samples required. Lastly, BAD’s explicit reliance on *common knowledge* limits the generality of the method.

In this paper we propose the *Simplified Action Decoder* (SAD), a method that achieves a similar goal to BAD, but addresses all of the issues mentioned above. At the core of SAD is a different approach towards resolving the conflict between exploration and being interpretable, which, like BAD, relies on the *centralized training with decentralized control* (CT/DC) regime. Under CT/DC information can be exchanged freely amongst all agents during *centralized training*, as long as the final policies are compatible with *decentralized execution*.

The key insight is that, during training we do not have to chose between being informative, by taking greedy actions, and exploring, by taking random actions. To be informative, the greedy actions do not need to be executed by the environment, but only need to be observed by the team mates. Thus in SAD each agent takes two different actions at each time step: One greedy action, which is not presented to the environment but observed by the team mates at the next time step as an additional input, and the “standard” (exploratory) action that gets executed by the environment and is observed by the team mates as part of the environment dynamics. Importantly, during greedy execution the observed environment action can be used instead of centralized information for the additional input, since now the agent has stopped exploring.

Furthermore, to ensure that these greedy actions and observations get decoded into a meaningful representation, we train an auxiliary task that predicts key hidden game properties from the action-observation trajectories. While we note that this method is in principle compatible with any kind of model-free deep RL method with minimal modifications to the core algorithm, we use a distributed version of recurrent DQN in order to improve sample efficiency, account for partial observability and reduce the risk of local optima. Similar to Value Decomposition Networks (VDN) (Sunehag et al., 2017) and QMIX (Rashid et al., 2018), we train a joint-action Q-function that consists of the sum of per-agent Q-values to allow for off-policy learning in this multi-agent setting.

Using SAD we establish a new state of the art for 2-5 players in Hanabi, with a method that not only requires less expert knowledge and compute, but is also more general than previous approaches. In order to ensure that our results can be easily verified and extended, we also evaluate our method on a proof-of-principle matrix game and plan to open-source our training code and agents. Beyond enabling more research into the self-play aspect of Hanabi, we believe these resources will provide a much needed starting point for the ad-hoc teamwork part of the Hanabi challenge.

2 RELATED WORK

Our work relates closely to research on emergent communication protocols using deep multi-agent RL, as first undertaken by Sukhbaatar et al. (2016) and Foerster et al. (2016). There has been a large number of follow-up papers in this area, so listing all relevant work is beyond the scope and we refer the reader to Nguyen et al. (2018), a recent survey on deep multi-agent RL. One major difference to our work is that the environments considered typically contain a cheap-talk channel, which can be modeled as a continuous variable during the course of training. This allows agents to, for example, use differentiation across the communication channel in order to learn protocols. In contrast, in our setting agents have to communicate through the observable environment actions themselves, requiring fundamentally different methods.

Furthermore, our work is an example of cooperative multi-agent learning in partially observable settings under centralized training and decentralized control. There have been a large number of papers in this space, with seminal work including MADDPG (Lowe et al., 2017) and COMA (Foerster et al., 2018), both of which are actor-critic methods that employ a centralized critic with decentralized actors. Again, we refer the reader to Nguyen et al. (2018) for a more comprehensive survey.

Until 2018, work on Hanabi had been focused on hand-coded methods and heuristics. Some relevant examples include SmartBot (O’Dwyer, 2019) and the so-called “hat-coding” strategies, as implemented by WTFWThat (Wu, 2018). These strategies use the information theoretic ideas that allow each hint to reveal information to all other agents at the same time. While they do not perform well for 2-player Hanabi due to the smaller action space, they get near perfect scores for 3-5 players.

In contrast, so far learning methods have seen limited success on Hanabi. Bard et al. (2019) undertake a systematic evaluation of current Deep RL methods for 2-5 players in two different regimes and open-source the Hanabi-Learning-Environment (HLE) to foster research on the game. They evaluate a feed-forward version of DQN trained on 100 million samples and a recurrent actor-critic agent with population based training using 20 billion samples. Notably both agents achieve near 0% win rate for 3-5 players in Hanabi. At a high level their DQN agent is a good starting point for our work. However, since the authors did not propose any specific method of accounting for the issues introduced by -greedy exploration in a ToM task, they resorted to setting it to zero after a short burn-in phase. The only state-of-the-art in Hanabi established by an RL agent is from Foerster et al. (2019) which we refer to in more detail in Section 1 and Section 4.

Recently there have also been attempts to train agents that are robust to different team-mates (Canaan et al., 2019) and even to extend to human-AI collaboration (Liang et al., 2019).

Poker is another partially observable multi-agent setting, although it is fundamentally different due to the game being zero-sum. Recent success in Poker has extensively benefited from search (Brown et al.). Examples of using search in Hanabi include Goodman (2019). For a more comprehensive review on previous results on Hanabi we refer the reader to Bard et al. (2019).

3 BACKGROUND

3.1 SETTING

In this paper we assume a Dec-POMDP (Oliehoek, 2012), in which agents interact in a partially observable environment. At each time step t , agent $i \in \{1, \dots, N\}$ obtains an observation $o_t^i = O(s_t; a)$, where $s_t \in S$ is the Markov state of the system and $O(s; a)$ is the observation function.

While our method are general, we restrict ourselves to turn based settings, in which at each time step only the acting agents takes an action, which is sampled from their policy $\pi^a(u^a; o^a)$, while all other agents take a no-op action. Here o^a is the action-observation history of agent $a = \{o_0^a; u_0^a; r_1; \dots; r_T; o_T^a\}$, T is the length of the episode and o^a are the weights of a function approximator that represents the policy, in our case recurrent neural networks, such as LSTMs (Hochreiter & Schmidhuber, 1997). Since we are interested in ToM, in our setting the observation function includes the last action of the acting agent, which is observed by all other agents at the next time step. We note that actions are commonly observable not only in board games but also in some real world multi-agent settings, such as autonomous driving.

As is typical in cooperative multi-agent RL, the goal of the agents is to maximize the total expected return, $J = E_{P(j)} R_0(\cdot)$, where $R_0(\cdot)$ is the return of the trajectory (in general $R(\cdot) = \sum_{t=0}^{\infty} \gamma^t r_t$) and γ is an optional discount factor. We have also assumed that agents are sharing parameters, as is common in cooperative MARL.

3.2 DISTRIBUTED RECURRENT DQN

In Q-learning the agent approximates the expected return for a given state action pair, assuming that the agent acts greedily with respect to the Q-function for all future time steps, $Q(s; u) = E_{P(j; s; u)} R_t(\cdot)$, where $u_t = u$ and $u_{t^0} = \arg\max_u Q(s_{t^0}; u)$; $t^0 > t$. A common exploration scheme is greedy, in which the agent takes a random action with probability ϵ and acts greedily otherwise. Importantly, the Q-function can be trained efficiently using the Bellman equation: $Q(s; u) = E_{s^0} r_{t+1} + \max_u Q(s^0; u)$, where for simplicity we have assumed a deterministic reward. In Deep Q-Learning (DQN) (Mnih et al., 2015) the Q-function is parameterized by a deep neural network and trained with experience replay.

In our work we also incorporate other best practice components of the last few years, including double-DQN (van Hasselt et al., 2015), dueling network architecture (Wang et al., 2015) and prioritized replay (Schaul et al., 2015). We also employ a distributed training architecture similar to the one proposed by Horgan et al. (2018) where a number of different actors with their own exploration rates collect experiences in parallel and feed them into a central replay buffer. Since our setting is partially observable the natural choice for the function approximator is a recurrent neural network. A combination these techniques was first explored by Kapturowski et al. (2019) in single agent environments such as Atari and DMLab-30.

3.3 CENTRALISED TRAINING, DECENTRALIZED EXECUTION AND JOINT Q-FUNCTIONS

The most straight forward application of Q-learning to multi-agent settings is Independent Q-Learning (IQL) (Tan, 1993) in which each agent keeps an independent estimate of the expected return, treating all other agents as part of the environment. One challenge with IQL is that the exploratory behavior of other agents is not corrected for via the \max operator in the bootstrap. Notably, IQL does typically not take any advantage of centralized training with decentralized control (CT/DC), a paradigm under which information can be exchanged freely amongst agents during the training phase as long as the policies rely only on local observations during execution.

There are various approaches for learning joint-Q-functions in the CT/DC regime. For example, Value-Decomposition-Networks (Sunehag et al., 2017) represent the joint-Q-function as a sum of per-agent contributions and QMIX (Rashid et al., 2018) learns a non-linear but monotonic combination of these contributions.

4 METHOD

4.1 THEORY OF MIND AND BAYESIAN REASONING

At the very core of interpreting the actions of another agent and ToM in general is Bayesian reasoning. Fundamentally, asking what a given action by another agent implies about the state of the world requires understanding why this action was taken. To illustrate this, we start out with an agent that has a given belief about the state of the world, given her current action-observation history, $B(s_t) = P(s_t | j^a)$.

Next the agent observes the action u_t^0 of her team mate and carries out a Bayesian update:

$$P(s_t | j^a; u_t^0) = \frac{P(u_t^0 | j^a; s_t) P(s_t | j^a)}{\sum_{s^0} P(u_t^0 | j^a; s^0) P(s^0 | j^a)} \quad (1)$$

$$= \frac{P(u_t^0 | j^a; s_t) B(s_t)}{\sum_{s^0} P(u_t^0 | j^a; s^0) B(s^0)}; \quad (2)$$

$$(3)$$

where for simplicity we have assumed that agent uses a feed-forward policy and observes the Markov state of the environment. Clearly, if an agent has access to the policy of their teammate during centralised training we could in principle evaluate the Bayesian belief explicitly. Instead, in this work we rely on RNNs in order to learn implicit representations of the sufficient statistics over the distribution of the Markov state given the action-observation histories.

4.2 EXPLORATION AND BELIEFS

Next we illustrate the impact of exploration on the beliefs, which we will do in the explicit (exact) case, since it serves as an upper bound on the accuracy of the implicit beliefs. Since we are looking at fully-cooperative settings we assume that the optimal policy of the agent is deterministic and any randomness is due to exploration. Given that we are focused on value based methods we furthermore assume an-greedy exploration scheme, noting that the same analysis can be extended to other methods. Under this exploration scheme $Q^a(u_t^a | s_t)$ becomes:

$$Q^a(u_t^a | s_t) = (1 - \alpha) I(u_t^a(s_t); u_t^{a^0}) + \alpha jU_j B(s_t); \tag{4}$$

where we have used $I(u_t^a(s_t))$ to indicate the greedy action of the agent $u_t^a(s_t) = \text{argmax}_u Q^a(u; s_t)$ and I is the indicator function.

While the first part corresponds to a filtering operator, in which the indicator function only attributes finite probability to those states that are consistent with the action taken under greedy execution, the exploration term adds a fixed (state independent) probability, which effectively 'blurs' the posterior:

$$P(s_t | u_t^a; u_t^{a^0}) = \frac{P(s_t) (1 - \alpha) I(u_t^a(s_t); u_t^{a^0}) + \alpha jU_j B(s_t)}{P(s^0) (1 - \alpha) I(u(s^0); u_t^{a^0}) + \alpha jU_j B(s^0)} \tag{5}$$

$$= \frac{(1 - \alpha) I(u_t^a(s_t); u_t^{a^0}) + \alpha jU_j B(s_t)}{jU_j + P(s^0) (1 - \alpha) I(u(s^0); u_t^{a^0}) B(s^0)} \tag{6}$$

$$= \frac{P(s_t) B(s_t)}{1 + jU_j \frac{P(s^0) (1 - \alpha) I(u(s^0); u_t^{a^0}) B(s^0)}{B(s^0)}} \tag{7}$$

$$+ \frac{(1 - \alpha) I(u_t^a(s_t); u_t^{a^0}) B(s_t)}{jU_j + P(s^0) (1 - \alpha) I(u(s^0); u_t^{a^0}) B(s^0)} \tag{8}$$

$$: \tag{9}$$

We find that the posterior includes an additional term of the form $B(s_t)$ which carries over an unfiltered density over the states from the prior. We further confirm that in the limit of $\alpha \rightarrow 1$, the posterior collapses to the prior $P(s_t | u_t^a; u_t^{a^0}) = B(s_t)$.

This can be particularly worrisome in the context of our training setup, whereby different agents run different, and potentially high throughout the course of training. It fundamentally makes the beliefs obtained less informative.

While not making the above argument explicitly, the Bayesian Action Decoder (BAD) (Foerster et al., 2019), resolves this issue by shifting exploration to the level of deterministic partial policies, rather than action-level, and tracking an approximate Bayesian belief. As outlined in Section 1 this comes at a huge cost in the complexity of the method, the compute requirements and in the loss of generality of the method.

4.3 SIMPLIFIED ACTION DECODING

In this paper we take a drastically simpler and different approach towards the issue. We note that the 'blurring', which makes decoding of an action challenging, is entirely due to the greedy exploration term. Furthermore, in order for another agent to do an implicit Bayesian update over an action taken, it is not required that this action is executed by the environment. Indeed, if we assume that other agents can observe the greedy action u_t^a , at every time step and condition their belief updated on this,

the terms depending on s_t disappear from the Bayesian update:

$$P(s_t | s_0; u) = \frac{P(s_t | s_{t-1}; u) B(s_t)}{P(s_0 | s_0; u) B(s_0)} \quad (10)$$

Therefore, to have our calculation, in the Simplified Action Decoder (SAD) the acting agent is allowed to 'take' two actions at any given time step. The first action, is the standard environment action, which gets executed as usual and is observed by all agents through the observation function at the next time step, as mentioned in Section 3. The second action is, the greedy action of the active agent. This action does not get executed by the environment but instead is presented as an additional input to the other agents at the next time step during training, taking advantage of the centralized training regime during which information can be exchanged freely.

Clearly we are not allowed to pass around extra information during centralized control but luckily this is not needed. Since we set s_0 at test time we can simply use the, now greedy, environment action obtained from the observation function as our greedy-action input.

Furthermore, to ensure that the agent meaningfully decodes the information contained in the greedy action, we also add an auxiliary supervised task to the training.

While this idea is compatible with any deep RL algorithm with minimal modifications, we use a recurrent version of DQN with distributed training, dueling networks and prioritized replay. We also learn a joint Q-function using VDN in order to address the challenges of multi-agent off-policy learning, please see Section 3 for details on all of these standard methods.

5 EXPERIMENTS

5.1 MATRIX GAME

We first verify the effectiveness of SAD in the two step, two player matrix game from Foerster et al. (2019), which replicates the communication challenge of Hanabi in a highly simplified setting. In this fully cooperative game each player obtains a privately observed 'card', which is drawn iid from two options (1,2). After observing her card, the first player takes one of three possible discrete actions (1, 2, 3). Crucially, the second player observes both her own private card and the team mate's action before acting herself, which establishes the opportunity to communicate. The payout is function of both the two private cards and the two actions taken by both agents, as shown in Figure 1. Importantly, there are some obvious strategies that do not require any communication. For example, if both player learn to play the 2nd action, the payout is always 8 points, independent of the cards dealt. However, if the players do learn to communicate it is possible to achieve 10 points on every pair of cards dealt.

Figure 1: Illustration of the matrix game from Foerster et al. (2019)

5.2 HANABI

Hanabi is a fully cooperative card game in which all players work together to complete piles of cards referred to as reworks. Each card has a rank, to 5, and a color G / B / W / Y / R. Each rework (one per color) starts with 5 cards and is finished once the 5 has been added. There are 5 of each rank and two of all other ranks for each of the colors, adding up to a total of 50 cards in the deck. The twist in Hanabi is that while players can observe the cards held by their team mates, they cannot observe their own cards and thus need to exchange information with each other in order to understand what cards can be played. There are two main means for doing so: First of all, players can take **gripped** actions in which they reveal the subset of a team mate's hand that matches a specific rank or color. An example hint is "Your third and fifth card are 5s". These hint actions cost **scarce** information tokens which can be replenished by **discarding** a card, an action that both removes the card from the game and makes it visible to all players.

Finally players can also choose to play a card. If this card is the next card for the rework of the corresponding color, it is added to the rework and the team scores one point. Otherwise the card is removed from the game, the identity is made public, and the team loses one of the 3 life tokens. If the team runs out of life tokens before the end of the game, all points collected so far are lost and the game finishes immediately. These rules result in a maximum score of 25 points in any game, which corresponds to all five reworks being completed with five cards per rework.

To ensure reproducibility and comparability of our results we use the Hanabi Learning Environment (HLE) (Bard et al., 2019) for all experimentation. For further details regarding Hanabi and the self-play part of the Hanabi challenge please see Bard et al. (2019).

5.3 ARCHITECTURE AND COMPUTE REQUIREMENTS

We borrow some ideas and insights from prior distributed Q-learning methods while bring innovations in our implementation to improve throughput and efficiency as well as extensions to MARL. Following Horgan et al. (2018) and Kapturowski et al. (2019), we use a distributed prioritized replay buffer shared by N asynchronous actors and a centralized trainer that samples mini-batches from the replay buffer to update the model. In each actor thread, we reinitialize environment sequentially and batch their observations together. The observation batch is then fed into an actor that utilizes a GPU to compute a batch of actions. All asynchronous actors share one GPU and the trainer uses another GPU for gradient computation and model updates. This is different from prior works which run single actor and single environment in each thread on a CPU. Our method enables us to run a very large number of simulations with moderate computation resources. In all Hanabi experiments, we run $N = 80$ actor threads with $K = 80$ environments in each thread on single machine with 40 CPU cores and 2 GPUs. Without this architectural improvement, it may require at least a few hundreds of CPU cores to run 6400 Hanabi environments with neural network agents and simulations have to be distributed across multiple machines, which will greatly hinder the reproducibility and accessibility of such research. Please refer to Appendix A for implementation details and hyper-parameters.

6 RESULTS

6.1 MATRIX GAME

As we can see in Figure 2, even in our simple matrix game the greedy action input makes a drastic difference. With an average reward of under 9 points, IQL does poorly in this task, just under-performing the Policy Gradient results from Foerster et al. (2019). In contrast, just by adding the greedy action as an additional input, we obtain an average performance of 9.50, matching the performance of BAD on this proof of principle task. Due to the simplicity of the task we are isolating just the impact of adding the greedy action to the input and are not investigating auxiliary tasks and recurrence in this setting. Results are averaged over 100 seeds, shading is error of the mean. The code is herbit.ly/2mBJLyk.

Figure 2: Results for the matrix game.

6.2 HANABI

As shown in Figure 3, our findings from the matrix game are for the most part confirmed on the challenging Hanabi benchmark. We compare SAD to three different ablations in order to illustrate the contributions of the different components. Baseline here is simply the recurrent DQN agent with parameter sharing. DQN is the same agent but also learns a joint Q-function and finally VDN+GreedyInputs the SAD agent without the auxiliary task. The main difficulty in evaluating the performance of SAD and the three ablations lies in the computational cost of running multiple seeds across 2-5 players, bearing in mind that each run uses around 8-10 billion samples and takes roughly 72 hours to run. Since we could only afford to run three seeds for each ablation for each number of players, we aggregate performance data across the different number of players in Figure 3.

Agent	2 Players		3 Players		4 Players		5 Players	
ACHA (Bard et al., 2019)	22.73	0.12	20.24	1.1%	21.57	0.12	16.80	0.13
	15.1%		1.1%		2.4%		0%	
BAD (Foerster et al., 2019)	23.92	0.01	-	-	-	-	-	-
	58.56%		-		-		-	
Best Seed: Baseline	23.89	0.01	23.59	0.01	22.59	0.01	20.64	0.01
	48.13%		39.04%		15.2%		1.73%	
Best Seed: VDN, SAD (w/ or w/o aux)	24.07	0.01	24.02	0.01	23.81	0.01	23.00	0.01
	56.41%		50.48%		38.53%		14.4%	

Table 1: Previous state of the art performances for learning methods in the unlimited regime of the HLE for 2-5 players compared to our best results and baseline.

This figure shows the average percentile that SAD and the three ablations obtain. The metric is normalized such that a method in which all three seeds are the lowest performing across all numbers of players will result in 0%, while one that consistently produces the best three seeds out of the 12 total seeds will achieve 100%. Further details on this metric are provided in Appendix B.

As we can see, under this metric SAD beats our baseline by a large margin, with clearly separated error bars. While the ablations indicate the effectiveness of the different components of the method, we note that additional runs are needed to test whether the auxiliary task contributes significantly to the performance. For completeness we have included a plot which shows all training runs for all numbers of players across all ablations and the method in the Appendix B.

Furthermore, as shown in Table 1, when we take the best achieved performance from our various training runs, we establish a new SOTA for learning methods on the self-play part of the challenge for 2-5 players, with the most drastic improvement being achieved for 3-5 players. In particular we beat both the ACHA agent from Bard et al. (2019) and the BAD agent, even though both of them use population based training and require more compute. We note that BAD was optimized for a different counting scheme, in which agents keep their scores when they run out of lives at the end of the game. This explains the higher win rate (58.6%) of BAD combined with a slightly lower mean score, when compared to SAD.

7 CONCLUSION AND FUTURE WORK

In this paper we presented the Simplified Action Decoder (SAD), a novel deep multi-agent RL algorithm that allows agents to learn communication protocols in settings where no cheap-talk channel is available. On the challenging benchmark Hanabi our work substantially improves the SOTA for an RL method for all numbers of players. For two players SAD establishes a new high-score across any method. Furthermore we accomplish all of this with a method that is both simpler and requires less compute than previous advances. While these are encouraging steps, there is clearly more work to do. In particular, there remains a large performance gap between the numbers achieved by SAD and the known performance of hand-coded strategies for 3-5 players. One possible reason is that SAD does not undertake any explicit exploration in the space of possible conventions. Another promising route for future work is to integrate search with RL, since this has produced SOTA results in a number of different domains including Poker, Go and backgammon.

REFERENCES

- Chris L Baker, Julian Jara-Ettinger, Rebecca Saxe, and Joshua B Tenenbaum. Rational quantitative attribution of beliefs, desires and percepts in human mentalizing. *Nature Human Behaviour* 1(4): 0064, 2017.
- 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. arXiv:1902.00506 [cs, stat], February 2019. URL <http://arxiv.org/abs/1902.00506> . arXiv: 1902.00506.
- Noam Brown and Tuomas Sandholm. Superhuman AI for heads-up no-limit poker: Libratus beats top professionals. *Science* pp. eaao1733, 2017.
- Noam Brown and Tuomas Sandholm. Superhuman AI for multiplayer poker. *Science* pp. eaay2400, 2019.
- Noam Brown, Tuomas Sandholm, and Brandon Amos. Depth-Limited Solving for Imperfect-Information Games. pp. 14.
- Murray Campbell, A Joseph Hoane Jr, and Feng-hsiung Hsu. Deep Blue. *Artificial intelligence*, 134 (1-2):57–83, 2002.
- Rodrigo Canaan, Julian Togelius, Andy Nealen, and Stefan Menzel. Diverse agents for ad-hoc cooperation in hanabi. arXiv preprint arXiv:1907.03840, 2019.
- Jakob Foerster, Ioannis Alexandros Assael, Nando de Freitas, and Shimon Whiteson. Learning to communicate with deep multi-agent reinforcement learning. *Advances in Neural Information Processing Systems* pp. 2137–2145, 2016.
- Jakob Foerster, Francis Song, Edward Hughes, Neil Burch, Iain Dunning, Shimon Whiteson, Matthew Botvinick, and Michael Bowling. Bayesian action decoder for deep multi-agent reinforcement learning. In *International Conference on Machine Learning* pp. 1942–1951, 2019.
- Jakob N Foerster, Gregory Farquhar, Triantafyllos Afouras, Nantas Nardelli, and Shimon Whiteson. Counterfactual multi-agent policy gradients. *Thirty-Second AAAI Conference on Artificial Intelligence* 2018.
- James Goodman. Re-determinizing information set monte carlo tree search in hanabi. arXiv preprint arXiv:1902.06075, 2019.
- Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Comput.* 9(8):1735–1780, November 1997. ISSN 0899-7667. doi: 10.1162/neco.1997.9.8.1735. URL <http://dx.doi.org/10.1162/neco.1997.9.8.1735> .
- Dan Horgan, John Quan, David Budden, Gabriel Barth-Maron, Matteo Hessel, Hado van Hasselt, and David Silver. Distributed prioritized experience replay. *CoRR* abs/1803.00933, 2018. URL <http://arxiv.org/abs/1803.00933> .
- Steven Kapturowski, Georg Ostrovski, John Quan, and Will Dabney. RECURRENT EXPERIENCE REPLAY IN DISTRIBUTED REINFORCEMENT LEARNING. pp. 1–19, 2019.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980, 2014.
- Claire Liang, Julia Proft, Erik Andersen, and Ross A Knepper. Implicit communication of actionable information in human-ai teams. *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems* pp. 95. ACM, 2019.
- Ryan Lowe, Yi Wu, Aviv Tamar, Jean Harb, OpenAI Pieter Abbeel, and Igor Mordatch. Multi-agent actor-critic for mixed cooperative-competitive environments. *Advances in Neural Information Processing Systems* pp. 6379–6390, 2017.

- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Belle-mare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharshan Kumaran, Daan Wierstra, Shane Legg, and Demis Hassabis. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 02 2015. URL <http://dx.doi.org/10.1038/nature14236>.
- Matej Moravčík, Martin Schmid, Neil Burch, Viliam Lisý, Dustin Morrill, Nolan Bard, Trevor Davis, Kevin Waugh, Michael Johanson, and Michael Bowling. Deepstack: Expert-level artificial intelligence in heads-up no-limit poker. *Science*, 356(6337):508–513, 2017.
- Thanh Thi Nguyen, Ngoc Duy Nguyen, and Saeid Nahavandi. Deep reinforcement learning for multi-agent systems: a review of challenges, solutions and applications. *arXiv preprint arXiv:1812.11794*, 2018.
- Arthur O’Dwyer. Hanabi. <https://github.com/Quuxplusone/Hanabi>, 2019.
- Frans A Oliehoek. Decentralized pomdps. In *Reinforcement Learning*, pp. 471–503. Springer, 2012.
- Tabish Rashid, Mikayel Samvelyan, Christian Schröder de Witt, Gregory Farquhar, Jakob N. Foerster, and Shimon Whiteson. QMIX: monotonic value function factorisation for deep multi-agent reinforcement learning. *CoRR*, abs/1803.11485, 2018. URL <http://arxiv.org/abs/1803.11485>.
- Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver. Prioritized Experience Replay. *arXiv e-prints*, art. arXiv:1511.05952, Nov 2015.
- David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016.
- David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go without human knowledge. *Nature*, 550(7676):354, 2017.
- David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dharshan Kumaran, Thore Graepel, et al. A general reinforcement learning algorithm that masters chess, shogi, and go through self-play. *Science*, 362(6419):1140–1144, 2018.
- Sainbayar Sukhbaatar, arthur szlam, and Rob Fergus. Learning Multiagent Communication with Backpropagation. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 29*, pp. 2244–2252. Curran Associates, Inc., 2016. URL <http://papers.nips.cc/paper/6398-learning-multi-agent-communication-with-backpropagation.pdf>.
- Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinícius Flores Zambaldi, Max Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z. Leibo, Karl Tuyls, and Thore Graepel. Value-decomposition networks for cooperative multi-agent learning. *CoRR*, abs/1706.05296, 2017. URL <http://arxiv.org/abs/1706.05296>.
- Richard S. Sutton. Learning to predict by the methods of temporal differences. *Mach. Learn.*, 3(1):9–44, August 1988. ISSN 0885-6125. doi: 10.1023/A:1022633531479. URL <https://doi.org/10.1023/A:1022633531479>.
- Ming Tan. Multi-agent reinforcement learning: Independent vs. cooperative agents. In *Proceedings of the tenth international conference on machine learning*, pp. 330–337, 1993.
- Hado van Hasselt, Arthur Guez, and David Silver. Deep reinforcement learning with double q-learning. *CoRR*, abs/1509.06461, 2015. URL <http://arxiv.org/abs/1509.06461>.
- Ziyu Wang, Nando de Freitas, and Marc Lanctot. Dueling network architectures for deep reinforcement learning. *CoRR*, abs/1511.06581, 2015. URL <http://arxiv.org/abs/1511.06581>.
- Jeff Wu. Hanabi simulation in rust. <https://github.com/WuTheFWasThat/hanabi.rs>, 2018.

A NETWORK ARCHITECTURE AND HYPER-PARAMETERS FOR HANABI

Our Hanabi agent uses dueling network architecture (Wang et al., 2015). The main body of the network consists of 1 fully connected layer of 512 units and 2 LSTM (Hochreiter & Schmidhuber, 1997) layers of 512 units, followed by two output heads for value and advantages respectively. The same network configuration is used across all Hanabi experiments. We take the default featurization of HLE and replace the card knowledge section with the V0-Belief proposed by Foerster et al. (2019). The maximum length of an episode is capped at 80 steps and the entire episode is stored in the replay buffer as one training sample. This avoids the “slate hidden states” problem as described in Kapturowski et al. (2019) so that we can simply initialize the hidden states of LSTM as zero during training. For exploration and experience prioritization, we follow the simple strategy as in Horgan et al. (2018) and Kapturowski et al. (2019). Each actor executes an ϵ -greedy policy where $\epsilon_i = \frac{1}{1 + \frac{1}{\epsilon_i}}$ for $i \in \{0, \dots, N-1\}$ but with a smaller $\epsilon = 0.1$ and $\epsilon = 7$. For simplicity, all players of a game use the same epsilon. The per time-step priority p_t is the TD error and per episode priority is computed following $p_e = \max_t p_t + (1 - \gamma)^t$ where $\gamma = 0.9$. Priority exponent is set to 0.9 and importance sampling exponent is set to 0.6. We use n -step return (Sutton, 1988) and double Q-learning (van Hasselt et al., 2015) for target computation during training. The discount factor γ is set to 0.999. The network is updated using Adam optimizer (Kingma & Ba, 2014) with learning rate $\alpha = 6.25 \times 10^{-5}$ and $\beta = 1.5 \times 10^{-5}$. Trainer sends its network weights to all actors every 10 updates and target network is synchronized with online network every 2500 updates. These hyper-parameters are fixed across all experiments.

In the baseline, we use Independent Q-Learning where each player estimates the Q value and selects action independently at each time-step. Note that all players need to operate on the observations in order to update their recurrent hidden states while only the current player has non-trivial legal moves and other players can only select ‘pass’. Each player then writes its own version of the episode into the prioritized replay buffer and they are sampled independently during training. The prioritized replay buffer contains 2^{17} (131072) episodes. We warm up the replay buffer with 10,000 episodes before training starts. Batch size during training is 128 for games of different numbers of players.

As mentioned in Section 4, the SAD agent is built on top of joint Q-function where the Q value is the sum of the individual Q value of all players given their own actions. One episode produces only one training sample with an extra dimension for the number of players. The replay buffer size is reduced to 2^{16} for 2-player and 3-player games and 2^{15} for 4-player and 5-player games. The batch sizes for 2-, 3-, 4-, 5-players are 64, 43, 32, 26 respectively to account for the fact that each sample contains more data.

Auxiliary task can be added to the agent to help it decode the greedy action more effectively. In Hanabi, the natural choice is the predict the card of player’s own hand. In our experiments, the auxiliary task is to predict the status of a card, which can be playable, discardable, or unknown. The loss is the average cross entropy loss per card and is simply added to the TD-error of reinforcement learning during training.

B MORE DETAILS ON EXPERIMENTAL RESULTS FOR HANABI

In this section we provide detailed experimental results for Hanabi. Table 2 is the complete version of Table 1 with ablations by removing auxiliary task, greedy action input and joint Q value one by one. Figure 4 shows the raw learning curves of different methods with 3 seeds per ablation / method. To obtain the curves shown in Figure 3, we first give each curve a score, ranging from 0 to 11, at every epoch based on its ranking among 12 curves (4 methods and 3 seeds each) in each graph. Then for each method, we average its ranking scores across all 12 curves (4 different numbers of players and 3 seeds each) and normalize its average score with 100% corresponding to the highest possible average score of 10 and 0% corresponding to the lowest possible average score of 1. These correspond to the extreme cases where one method produces the lowest or highest results for every single seed across all numbers of players.

Agent	2 Players		3 Players		4 Players		5 Players	
SmartBot (O’Dwyer, 2019)	22.99	0.00	23.12	0.00	22.19	0.00	20.25	0.00
	29.6%		13.8%		2.08%		0.0043%	
WTFWThat (Wu, 2018)	19.45	0.03	24.20	0.01	24.83	0.01	24.89	0.00
	0.28%		49.1%		87.2%		91.5%	
Rainbow (Bard et al., 2019)	20.64	0.03	18.71	0.01	18.00	0.17	15.26	0.18
	2.5%		0.2%		0%		0%	
ACHA (Bard et al., 2019)	22.73	0.12	20.24	1.1%	21.57	0.12	16.80	0.13
	15.1%		1.1%		2.4%		0%	
BAD (Foerster et al., 2019)	23.92	0.01	-	-	-	-	-	-
	58.56%		-		-		-	
Baseline	23.89	0.01	23.59	0.01	22.59	0.01	20.64	0.01
	48.13%		39.04%		15.2%		1.73%	
VDN	23.95	0.01	24.02	0.01	23.66	0.01	21.20	0.01
	49.17%		50.48%		32.51%		2.33%	
VDN with GreedyInput	23.96	0.01	23.84	0.01	23.81	0.01	23.00	0.01
	48.22%		44.33%		38.53%		14.4%	
SAD	24.07	0.01	23.78	0.01	23.49	0.01	21.53	0.01
	56.41%		48.46%		34.09%		2.39%	

Table 2: Performance of various methods on Hanabi, evaluating the best of 3 seeds for each of our method and ablations. Each model is evaluated on 100K games with different seeds. The s.e.m. is less than 0.01 for most models. Bold numbers are the best results achieved with learning algorithms.

