Generating and Adapting to Diverse Ad-Hoc Cooperation Agents in Hanabi

Rodrigo Canaan NYU New York, USA rodrigo.canaan@nyu.edu Xianbo Gao NYU New York, USA xg656@nyu.edu

Julian Togelius

NYU

New York, USA
julian.togelius@nyu.edu

Andy Nealen *USC*Los Angeles, USA nealen@usc.edu

Stefan Menzel

HRI Europe GmbH

Offenbach, Germany
stefan.menzel@honda-ri.de

Abstract—Hanabi is a cooperative game that brings the problem of modeling other players to the forefront. In this game, coordinated groups of players can leverage pre-established conventions to great effect, but playing in an ad-hoc setting requires agents to adapt to its partner's strategies with no previous coordination. Evaluating an agent in this setting requires a diverse population of potential partners, but so far, the behavioral diversity of agents has not been considered in a systematic way. This paper proposes Quality Diversity algorithms as a promising class of algorithms to generate diverse populations for this purpose, and generates a population of diverse Hanabi agents using MAP-Elites. We also postulate that agents can benefit from a diverse population during training and implement a simple "meta-strategy" for adapting to an agent's perceived behavioral niche. We show this meta-strategy can work better than generalist strategies even outside the population it was trained with if its partner's behavioral niche can be correctly inferred, but in practice a partner's behavior depends and interferes with the meta-agent's own behavior, suggesting an avenue for future research in characterizing another agent's behavior during gameplay.

I. INTRODUCTION

This paper addresses the problem of generating agents for *ad-hoc cooperation* [1] in games. In ad-hoc cooperation, an agent must be able to achieve high performance in a cooperative task when paired with one or more partners, called *ad-hoc teammates*, whose identity is not known prior to the start of the task, and with no prior coordination of strategies. This type of problem can model many real-world interactions between human actors and autonomous artificial agents, such as a driver navigating a busy street among other vehicles and pedestrians or a digital assistant giving personalized suggestions to a human user.

This precludes the agent from taking advantage of assumptions and conventions that could be taken for granted in scenarios where one's partners are known in advance (which includes self-play). On the other hand, online modeling of other agents takes greater importance. A typical challenge in agent modeling is about *predicting* the future actions of other agents, but in environments with asymmetrical information, it also becomes important to *interpret* observed actions and infer what they might imply about hidden features of the world. In essence, our agents should be able to represent the distinct mental state of other actors and see the world from their perspective. This ability has also been referred to as

having a theory of mind [2]. Hanabi is a game that brings these issues to the forefront due to the unique nature of its asymmetrical information, where players don't see their own cards and communication is heavily restricted. We describe the game in section II.

Evaluation of agents in an ad-hoc setting typically occurs by sampling teammates from a pool of potential partners. However, the choice of which agent to include in this pool, how to account for the behavioral diversity of these agents and whether the pool is made public or secret to the designers of ad-hoc agents has important implications to the generality and reproducibility of a benchmark or competition. We discuss these issues, along with other related work, in section III.

Sections IV and V contain the methodology and results of the first main contribution of this paper: metrics for characterizing the behavior space of Hanabi agents and a method for generating diverse, high-quality agents using MAP-Elites [3], which attempts to fill out each niche of this space with rule-based agents having as high fitness (given by self-play performance) as possible.

The remainder of the paper contains our second main contribution: designing an ad-hoc cooperation "meta-agent" that leverages the pools of partners generated by MAP-Elites. We do this by pre-computing a generalist strategy for the pool as a whole and a specialized strategy for each partner in the pool, then attempting to estimate a partner's niche during gameplay and selecting the action suggested by the appropriate specialized strategy. We show that this adaptation strategy can perform better than the generalist strategy if a partner's niche can be correctly estimated, but in practice the meta-agent is unable to accurately estimate the niche, due to what seems to be an issue of interference between the meta-agent's behavior and its partner's behavior.

This paper is an extension of a previous paper [4] presented at the 2019 IEEE Conference on Games (CoG). Sections VI and VII are novel contributions describing the meta-agent and associated results. The earlier sections include, among smaller changes, the addition of pseudocode describing the MAP-Elites algorithm, the definition of a new behavioral metric (IPP) that doesn't rely on secret information and a more detailed discussion about the relevance of techniques that account for behavioral diversity in the context of ad-hoc cooperative game AI benchmarks.

II. HANABI: THE GAME

Hanabi is a cooperative card game designed by Antoine Bauza and has won the prestigious *Spiel des Jahres* award for tabletop games in 2013. It is played by groups of 2-5 players who try to play stacks of cards in correct order of rank or value (from 1 to 5) for each of the five colors in the game (B, R, Y, W and G). Players play with the contents of their hands facing outwards, so that each player sees the cards every other player has, but not their own cards. The group can only communicate through hint actions, which allow the current player to select another player and point to all cards of a chosen rank or color in their hand, at the expense of a hint token from a shared pool. When a card is played to the table, the group scores a point if it is the next card in its color stack, or loses a life otherwise. A player can also discard a card from their hand, which recovers one information token.

The game ends with a victory for the group if all five stacks are completed with cards ranked 1 to 5 of that color. Whenever a card is played or discarded, players must draw back to their hand limit from a draw deck. If the draw deck is exhausted, every player gets one last turn to take an action, after which the game ends in defeat if not all stacks are complete. The game also ends with defeat if the group loses three lives. The score of the game is the number of cards successfully played by the group, that is, 25 in case of victory or from 0 to 24 in case of defeat.

Hints provide grounded information by disclosing the rank or color of cards. Examples of hints are "your first, second and fourth card are 1's", "your middle card is a 5" and "your two rightmost cards are yellow". But in addition to this grounded layer of meaning, players try to infer additional implicit meaning from each hint by taking into account a model of the other player. For example, if the 1's of some (but not all) colors have been played, most players who receive a hint of "your leftmost card is a two" would assume that it refers to a playable card, even though nothing was said of its color. The player giving the hint should therefore *predict* how the receiving player would act in different scenarios and the receiving player must in turn *interpret* what each hint (and other actions) says about the state of the game.

Over time, conventions can either emerge organically or be formally agreed to in a group. Conventions provide a guideline for how hints ought to be interpreted and which hints should be given in each situation. Examples of conventions are "hints should, if possible, identify playable cards rather than unplayable ones" and "players should discard unidentified cards from oldest to newest". But these can easily backfire if all players are not using the same convention. For AI research, this makes the self-play problem, where all agents are known to be following the same strategies and conventions, fundamentally different from the problem of ad-hoc cooperation.

III. RELATED WORK

A. Hanabi-Playing agents and the CIG/CoG competition

Many of the early approaches [5]-[7] for playing Hanabi with AI were variations of a simple strategy for self-play

which prioritizes playing cards that are believed to be playable, followed by giving hints that identify playable cards in other player's hands, followed by discarding cards that are believed to not be necessary. Walton-Rivers *et al.* [8] re-implemented many of these agents under a rule-based paradigm, while addressing the problem of playing with a diverse population of agents with different strategies in Hanabi.

A rule-based agent is defined by an ordered sequence of rules. Each rule takes a game state and checks whether a condition is true. If the condition is true, the rule also specifies which action to take. Otherwise, the agent moves on to the next rule in order. Rules can be divided into play rules (e.g. "play a card estimated to be playable with probability > 60%"), tell rules (e.g. "give new information about a playable card") and discard rules (e.g. "discard the oldest card in hand"). Their evaluation of an agent was based on that agent's average score across all parings with agents from a fixed pool.

The 2019 Hanabi CoG competition [9] is based on this work. It first took place at the 2018 CIG conference, where participants submitted agents both for a self-play (or "mirror") track and a mixed play track. In the mixed play track, the agents had to play with a pool similar to the one used in [8], but the exact agents were not made public before the competition. The winner of the competition was a variant of Monte Carlo Tree Search [10] by Goodman [11] designed to deal with problems of strategy fusion and nonlocality that arise from executing tree search in a hidden information environment. They also use neural networks to model a pool of other sample agents, and bayesian updates keep track of which agent in this pool best approximates the current partner. It achieved a score of 13.28 in the mixed track and 20.57 in the mirror track [12].

Our 2018 competition entry, which took second place, is described in [13]. We implemented an evolutionary algorithm to make rule-based agents by searching for a well-performing sequence of rules both for self-play and mixed play, using the same pool as [8] for mixed play. It achieved a score of 12.85 in the mixed track and 17.52 in the mirror track. While the current paper is based on that previous work, here we ignore the standard pool used in [8] and [13] and procedurally create our own pool, whose agents are evaluated by how well they fare on self-play and when paired with each other.

The 2019 competition also featured a mirror and mixed track, plus a new learning track where agents can adapt to repeated play with the same partners. While this paper does not aim directly to develop agents to compete, we hope the pools of ad-hoc partners we are generating can help in evaluating future agents, and in better understanding what types of play work better with agents exhibiting a variety of behaviors.

Other than the CoG competition and its related agents, another body of work on Hanabi agents focuses on Reinforcement Learning (RL) agents, starting with the *Hanabi Learning Environment* by Bard *et al.* [14]. Their Actor Critic Hanabi Agent (ACHA) introduced in the paper scores over 20 points in self-play, but scores close to zero when paired with independently trained instances of itself and with a Rainbow

DQN [15] agent. They also propose an ad-hoc setting where self-play playtraces of the partner agent are provided prior to gameplay for learning, but to our knowledge there are no published agents taking advantage of this feature.

The best-performing self-play agents we are aware of is by Lerer *et al.* [16]. This is a hybrid agent that combines a pre-established "blueprint policy" with a public search procedure. The best version of the agent uses an RL agent called the Simplified Action Decoder (SAD) [17] as blueprint policy. SAD achieves an average score in 2-player games of 24.08, and Lerer's agent using SAD as blueprint policy achieves an average score of 24.61.

Both SAD and Lerer's search procedure benefit greatly from coordination taking place before the game, and as such are not expected to perform well in ad-hoc scenarios and have been noted to learn conventions that arbitrarily map color hints to playing or discarding a card in specific positions. Conventions of this type had been previously used in agents based on the "hat-guessing" paradigm first proposed by Cox *et al.* [18] and expanded on by Bouzy [19] and Wu [20], which also achieve scores over 24 in games with 3 or more players.

Studies involving evaluation with human players in Hanabi include the rule-based agents by Eger *et al.* [7] and Liang *et al.* [21] inspired by Grice's maxims [22] and a recent method for dealing with symmetries of a Dec-POMDP called Other Play [23], which used SAD as the underlying architecture for Hanabi play. All of these studies achieve similar mean scores around 15 with humans.

B. Ad-Hoc Cooperation

Ad-hoc Cooperation is the problem of cooperating with arbitrary partners, with no prior coordination. We call those partners ad-hoc teammates (human or AI agents). One of the earliest formalizations of this problem involving teams of autonomous artificial agents is by Stone et al. [1].

Their framing of the problem can be stated as such: given a pool of potential teammates A and a domain of potential tasks D, create an agent a that maximizes the expected payoff when asked to perform a randomly sampled task $d \in D$, paired with a randomly sampled subset of teammates $B \subset A$. Importantly, B is not known at the start of each individual task. In our case, the domain D is the set of starting configurations of a game of 2-player Hanabi, which includes randomizing the choice of starting player.

An important difference between this formulation and the Hanabi challenges proposed in [9] and [14] is that in Stone *et al.* [1], the pool of potential agents A is assumed to be known in advance (although it could possibly be infinite), whereas in [9] and [14] this pool is secret.

Comparing these two formulations, treating A as given makes the ad-hoc performance of an agent easier to reproduce and to compare with other agents by simply making a large enough sample of teammates from A and tasks from D, but high performance in this settings might come as a result of over-specialization to the partners in A and might not generalize to other choices of partners.

On the other hand, treating A as unknown introduces problems for reproducibility and comparison of performance with other agents. How could we compare a new agent's performance to those of, say the participants of the competition described in [9] after the paired controllers used in that competition have been made public? How could we compare experiments performed with different choices of A? How could we isolate an agent's intrinsic adaptation abilities from its designer's prior beliefs about the unknown agents in A?

In either formulation of the problem, the designer of the benchmark ideally wants to avoid a situation where a single non-adaptive agent plays well with all partners in the pool. Moreover, if the benchmark is meant to serve as a proxy for cooperation with humans, it requires some characterization of the space of human behaviors and a way to make sure A contains agents representative of this space.

A main contribution of this paper is the proposal of a method for generating the pool of partners A while explicitly taking into account both the behavioral diversity and the quality of the evaluation pool, which addresses many of these issues. We do this through the use of Quality Diversity algorithms such as MAP-Elites, which we discuss next.

C. Quality Diversity and MAP-Elites

Quality Diversity [24] (QD) algorithms are a class of population-based search algorithms that aim to generate a large number of solutions that are behaviorally diverse and of high quality. Behaviorally diverse means the agents are distributed representatively across a behavior space induced by one or more behavioral metrics. High quality means the agent performs well according to some fitness function.

Diversity of behavior can be pursued either as desirable target in its own right or as an intermediate step to high-quality solutions in deceptive fitness landscapes, as showcased by novelty search [25]. QD differs from novelty search, however, because it does not optimize for novelty alone, but searches for both behavioral diversity and high fitness at once. QD also differs from Multi-Objective Optimization [26], which searches for trade-offs between one or more objectives, because QD actively attempts to find high-quality solutions in all regions of the behavior space, not just those with good trade-offs.

MAP-Elites [3] is an example of QD algorithm that attempts to "illuminate" the behavior space by mapping each individual to a behavioral "niche", while maintaining an archive of the best individual (an elite) in each niche. MAP-Elites was first proposed to pre-compute a variety of effective gaits for a six-legged robot, so that, when the robot suffers damage, it can quickly search for a gait that adapts to the damage and allows it to keep moving at a decent pace [27].

IV. GENERATING DIVERSE AGENTS WITH MAP-ELITES

Our approach starts with using a QD algorithm such as MAP-Elites to generate a pool of high-quality, behaviorally diverse agents. This pool is used as our choice of A in the

formulation described in section III-B and addresses the issues raised in that section in the following ways:

- The creator of the pool needs to specify only the behavioral metrics that characterize a solution and the optimization method, but not the individual teammates that go in the pool, potentially reducing human bias in the selection of these teammates.
- Agents are guaranteed to be diverse with respect to the chosen behavioral metrics, which minimizes the risk that a single strategy cooperates well with all of them.
- The choice of behavioral metrics and its induced behavior space can also help with player modelling: we can ask where humans fall in a given space or which metrics best capture differences between individual human players.
- Generating the evaluation pool stochastically with a given distribution provides a middle ground between the setting where the complete pool is known in advance (which might encourage over-specialization) and where the pool is secret (which brings issues of reproducibility).

We believe our approach is useful both for the individual researcher, who can use it to create a pool of training partners in scenarios where the evaluation pool is secret, and to the larger community when designing new benchmarks and competitions that are meant to test the ability to cooperate with a variety of strategies or serve as a proxy for cooperation with humans.

Below, we describe how MAP-Elites was implemented to this end. The implementation is available in our public github repository ¹.

A. Definition of the feature space

MAP-Elites requires us to select one or more *behavior metrics* that induce a *behavior space* and serve as coordinates to locate any agent in that space.

We then need to *discretize* this space, which defines *niches* where all agents have all behavior metrics within a certain range.

Finally, we use some optimization protocol, such as evolutionary search, to simultaneously search for individuals with as high fitness as possible within each niche. This is done by keeping an archive of *elites* of each niche, representing the best agent found so far within that niche, then comparing a newly evaluated candidate's fitness with the fitness of the elite in the same niche as the candidate.

We chose the following behavioral metrics as dimensions of the behavior space in Hanabi:

• Information per Play (IPP): whenever an agent plays a card, we verify whether it knows its color and/or rank. Each of these is considered one piece of information. We count how many pieces of information the agent knows (either 0, 1 or 2) for each card that is played, then average this value across all played cards. Finally, we divide by 2 to get a number between 0 and 1. An agent scoring 1 in this dimension only plays cards that are fully known (both

- color and rank), whereas an agent scoring zero would only play cards it knows nothing about.
- Communicativeness: defined as the fraction of time an agent will choose to give a hint if a hint token available at the start of the turn. An agent scoring 1 in this dimension would always give a hint if possible, being fully communicative, while an agent scoring 0 would never give any hints.

In [4], we used a metric called Risk Aversion instead of IPP. It reflected the average probability that a card is playable, from the perspective of the agent, across all played cards. However, this probability depends on which cards an agent sees in its partners' hands, and this information is not available from the perspective of each partner. This makes Risk Aversion a hard metric to estimate for other players with the information available during gameplay. For this reason, we decided to replace it with IPP, which depends only on public information.

These dimensions were chosen because they are easy to measure and we believe that they are strategically meaningful, requiring different strategies to play with at different points in the feature space. In particular, both Risk Aversion and IPP were meant as proxies for a pattern that can be observed in game-play between humans: inexperienced players usually only play cards they know everything about (thus having values of Risk Averion or IPP close to 1), where more experienced players are more often comfortable playing cards under partial information as a result of either accounting for other player's apparent beliefs and intentions (theory of mind) or of following a particular pre-established convention.

IPP has the added benefit over Risk Aversion that well-performing agents with low IPP are theoretically possible, if the agent is using a convention that often allows a card that was never pointed at to be played (such as "if I give a hint of the color Yellow, you should play your second card"), where a low value of Risk Aversion forces the agent to only play cards that are known to be very unlikely to be playable.

We suspected that the highest-scoring behavior in self-play would fall at some value much greater than 0, but lower than 1 for both dimensions: a 0 in either dimension leads to obviously degenerate play, but good play likely requires playing cards under some uncertainty (implying IPP and risk aversion < 1) and sometimes passing up the opportunity to give a hint so that the other player can better utilize the hint token (implying Communicativeness < 1).

Note also that while these dimensions help describe an agent's play, they don't completely determine it. Communicativeness does not tell us *which* hint will be given, only the likelihood that *some* hint will be given if a hint token is available. Similarly, IPP does not tell us *whether* the agent will play a card, only how much is known on average about it given that it was played.

Each metric takes values in the range of [0,1], and we chose to discretize them at intervals of 0.05, defining 20 intervals in each dimension of the behavior space, which amounts to a total of 400 niches.

¹https://github.com/rocanaan/Hanabi-Map-Elites

B. Representation and operators

We use a similar representation of individuals as the one we used in [13]. Each individual is represented by a chromosome defined by a sequence of 15 integers, each integer representing one of 135 possible rules. An agent's action is determined by simply moving through the rules in the order they appear in the chromosome and selecting the action returned by the first rule that applies. An agent might have rules that never fire during gameplay (for example, a rule that says "discard a random card" would never fire if it comes after "discard your oldest card"). An agent can also have duplicate rules, in which case the second instance of the rule will never fire (assuming the rule either fires or not deterministically, which is true for the rules we are using). Nevertheless, these unused or repeated rules are part of an agent's genetic representation and can be passed on to its offspring. We selected 15 as chromosome length because our agents from [13] rarely had more than 10 different rules activated.

The first few chromosomes (in our experiments, 10⁴) are implemented by sampling rules uniformly at random from the ruleset, while the remaining chromosomes are generated by mutation and crossover of the elite in a random niche. Mutation is implemented by randomly replacing each rule in a chromosome with a random rule with probability 0.1. Crossover happens with probability 0.5 and is implemented by selecting another individual from the population and randomly selecting (with probability 0.5) the corresponding rule from either parent at each gene.

C. Pseudocode of the MAP-Elites algorithm

With these metrics, representation and operators in mind, algorithm 1 shows the abstracted pseudocode of the MAP-Elites algorithm.

Algorithm 1: MAP-Elites

```
 \begin{aligned} & \textbf{Result:} \ A, \text{ an archive with each niche's elite.} \\ & generation \leftarrow 0; \\ & A \leftarrow \emptyset; \\ & \textbf{while } generation < G \textbf{ do} \\ & c \leftarrow newChromossome(A, generation) ; \\ & x \leftarrow makeAgent(c); \\ & f \leftarrow fitness(x); \\ & i, j \leftarrow niches(x); \\ & \textbf{if } A_{i,j} = null \textbf{ then} \\ & | A_{i,j} \leftarrow c; \\ & \textbf{else} \\ & & | f_{elite} \leftarrow fitness(makeAgent(E_{i,j})) ; \\ & \textbf{if } f > f_{elite} \textbf{ then} \\ & | A_{i,j} \leftarrow c \end{aligned}
```

In our experiments, the functions *newChromosome*, *makeAgent*, *niches* are implemented as described below:

If generation $< 10^4$, newChromosome returns a list of 15 rules by sampling the rule-set uniformly. Otherwise, the

new chromosome is generated by mutation and crossover of random parents sampled from A, as described in section IV-B.

makeAgent simply returns an agent object that follows a policy determined by applying the chromosome rules in order, as also described in section IV-B.

fitness returns the average score of the agent after playing 100 matches in self-play mode. While these matches are played, a number of statistics can be recorded, such as the number of hints given, the number of turns where a hint token was available, the total number of cards played and how many pieces of information was known about each played card.

niches calculates the Communicativeness and IPP values based on these stored statistics, and converts these behavior metrics (which take values between 0 and 1) into two integer indexes, according to how many niches each metric's range is divided into. In our case, we divided each metric in 20 equally-sized parts, so a Communicativeness value between 0 and 0.05 corresponds to the first niche on the Communicativeness dimension, a value between 0.05 and 0.1 to the second and so on, until the last niche corresponding to values between 0.95 and 1.

After that, the program checks whether an elite has already been assigned to the corresponding entry $A_{i,j}$. If that entry is empty, the chromosome of the current agent is stored in that cell. Otherwise, we re-calculate the fitness of the elite in $A_{i,j}$ by instantiating its corresponding agent and computing its average score over 100 games, with the same random seed as used to play the games involved.

We do this rather than simply storing the elite's fitness value because an agent's fitness is obtained through a stochastic process, so it would be possible to overestimate an elite's fitness if it simply had an unusually lucky set of games, making it too hard for future generations to replace.

V. MAP-ELITES RESULTS

A. Self-play Scores

We used Map-Elites to generate and evaluate agents by executing three separate runs of algorithm 1. Each run generated and evaluated a total of 10^6 candidate individuals and recorded the chromosome of the elite in each of the 400 behavioral niches where a score greater than zero was achieved. A run's coverage is defined as the number of niches successfully filled by an elite in the run.

While during evolution each individual's fitness was estimated by playing one hundred games, after evolution we reevaluated each elite by playing a thousand self-play games. Table I shows the maximum self-play score of any elite in each run during this re-evaluation, the average self-play score of agents in all covered niches and each run's coverage.

All three runs had very similar coverage of either 310 or 311 out of 400. The mean score over covered niches was also very similar from run to run, with an average of 10.04 across the three runs. The score of the best agent in each run varied from 19.99 to 20.43. Across all 932 agents, the average Standard Deviation (SD) of their score was 2.71 and the maximum SD of the score of any agent was 7.15, corresponding to a

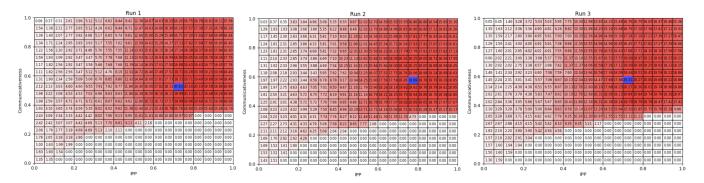


Fig. 1. Main results of the MAP-Elites experiment after reevaluating each elite in each run for 1000 games each. Values represent the fitness (score) of the best individual in that niche, with redder entries corresponding to higher scores. The maximum score for each run is highlighted in blue.

Run	# Candidates	Max Score	Average Score	Coverage
Run 1	1000000	20.43	10.02	311
Run 2	1000000	19.99	9.99	311
Run 3	1000000	20.31	10.11	310

TABLE I

Number of candidates generated, self-play score of the best performing elite, average score (over covered niches) and coverage of each of the three runs. The highest s.e.m. of any of the 932 agent's score is 0.23.

maximum Standard Error of the Mean (S.E.M.) of 0.23 with 1000 matches played.

Figure 1 shows the fitness of the elite in all 400 niches of each run. The three runs showed a very similar fitness landscape, with highest scores concentrated in the region with high values of both Communicativeness and IPP (upper right of each graph). The region with high IPP, but low Communicativeness were largely not covered by any agent scoring over 0 in self-play (bottom right of the graphs). An agent in this region would rarely give hints, yet only play cards it knows a lot of information about. In self-play, such agent would never give enough hints to its partner to satisfy their high IPP requirement and would thus score zero.

Note that, while the well-performing agents have high values for both behavior metrics, the best agent in each run was not found at the extreme of either metric. In run 1, the best agent was at the niche centered around (0.725,0.525), corresponding to IPP between 0.7 and 0.75 and Communicativeness between 0.5 and 0.55. This agent had an average self-play score of 20.43. For runs 2 and 3, the niche of the best agent was centered around (0.775,0.575) and (0.725,0.575), with scores of 19.99 and 20.33 respectively. The score of the best agent in runs 1 and 3 are an improvement over our best 2-player agent (called Mirror Situational) reported in [13], which had a reported self-play score of 20.07.

B. Diversity between runs

One of our goals with this work is to address the issues of reproducibility related to the choice of the pool A as described in section III-B by generating this pool in a stochastic way, offering a middle-ground between a public pool and a secret one. This requires our agents to not only be diverse within the

population generated by a single run, but also that the agents be reasonably diverse from run to run.

We evaluated this by comparing what we call "corresponding pairs" from the three runs. These are the pairs of elites that occupy the same niche in any two different runs. We performed experiments measuring how well these corresponding agents play together, how similar their chromosomes are and how similar the action they take are (given the same game-state).

For the first measurement, we simulated 1000 2-player games between each corresponding pair. There are 311 corresponding pairs between runs 1 and 2, 310 pairs between runs 1 and 3 and 310 pairs between runs 2 and 3 (due to the slightly smaller coverage of run 3). Table II shows the results of this experiment. For each niche, we call the score of a corresponding pair between two runs that pair's crossplay score. On average, the cross-play score was 10.02 across all pairs in the three runs. It is useful to compare that number with the average self-play score of agents in all the 3 runs (10.04). This means pairs of corresponding agents performed, on average, about as well with each other as they perform on self-play.

This might be an indication that these agents are in fact too similar from run to run. To test that, we compared their chromosomes directly. For each pair, we computed the Hamming distance [28] between each agent's chromosomes. Each chromosome was represented by 15 genes, each gene mapping to one of 135 possible rules. We computed the average Hamming distance to be 14.26, meaning each pair of chromosomes shared, on average, fewer than 1 rule in the same position.

While the chromosomes are very different, it would still be possible that the agents exhibited very similar behavior (phenotype) despite the distance between their genotypes.

To account for that, we investigated how often two agents from different populations that occupy the same niche take the same actions, if given the same game state. We did this by first playing 10 self-play games with each of the agents of population 1 and recording each game state, for a total of 149611 game states. Then, for each of these game states, and for each of the 310 niches where all runs of the MAP-Elites algorithm produced a valid elite, we verified whether

Run	1	2	3
1	10.02	9.95	10.02
2	-	9.99	10.10
3	-	-	10.11
	TAB		

CROSS-PLAY SCORE OF PAIRS OF CORRESPONDING AGENTS IN EACH RUN.
THE DIAGONAL IS THE SAME AS THE SELF-PLAY SCORES REPORTED IN
TABLE I. THE HIGHEST S.E.M. OF THE ANY OF THE 931 INDIVIDUAL
PAIRINGS IS 0.23

the two corresponding agents from two different populations take the same action. On average, these corresponding pairs took the same action on 62% of the game states we analysed. Considering there are up to 20 legal actions per game state (and an average of 13 legal actions), this is a much higher similarity than one would expect between two completely uncorrelated agents, but not enough to perfectly predict an agent's next play based solely on their niche.

In our previous paper [4] we also measured the similarity of the chromosomes and actions of agents within each run. We found that agents close to each other on the behavior map (figure 1) had similar chromosomes and chose similar actions, but agents further away from each other were more dissimilar. For space considerations, we omit these experiments here.

VI. THE META-AGENT

Next, we want to leverage the populations of diverse agents we generated to create more general agents for Ad-Hoc teamplay. In the Future Work section of our previous paper [4], we suggested an approach for building an adaptive agent based on identifying its partner's behavioral metrics and selecting an appropriate pre-computed strategy based on that behavior. We call this approach a meta-agent, since it attempts to select, during gameplay, a policy it estimates to work well with the current partner. We now describe this approach in more detail.

A. Offline training

During offline training, the agent must compute:

- A policy that maximizes the score with randomly sampled partners from the pool. We call this the generalist policy, and the meta-agent follows this policy if it is too uncertain about the behavioral metrics of its partner.
- For each region of the behavior space, a policy that plays well with partners from that region. We call each of these a **specialized policy**, to be used when the meta-agent estimates its partner's behavioral metrics to fall within that region.

Note that in principle, we could use any algorithm to generate a generalist or specialized policy. In our initial experiments, however, we select these policies from among the populations generated by MAP-Elites. We do this by first running 300 games between each pair of agents in a population and recording the corresponding average score. We call each pair of agents a match-up.

Then, we choose the generalist policy to be the policy of the agent that achieves the best average score across all of its match-ups, and the specialized policy for each niche as the policy of the agent who achieves the best score when paired with the agent in that niche.

More formally, let A be a (two-dimensional) population of agents. If we let $score(A_{i,j},A_{m,n})$ represent the expected score of the match-up between agents in positions (i,j) and (m,n) in A, we can define an agent's intra-population score to be:

$$intra(A_{i,j}) = \frac{1}{||A||} \sum_{m,n} score(A_{i,j}, A_{m,n})$$
 (1)

Where ||A|| is the number of valid agents in A (that population's coverage). Then, if we denote by (i_G, j_G) the indexes of the generalist policy, it follows that

$$(i_G, j_G) = \underset{i,j}{argmax}(intra(A_{i,j}))$$
 (2)

And the specialized policy (or specialized response) for niche (m,n) is given by the agent with indexes

$$(i_{m,n}, j_{m,n}) = \underset{i,j}{\operatorname{argmax}}(\operatorname{score}(A_{i,j}, A_{m,n}))$$
 (3)

We call the collection of specialized policies and their indexes a **response table**, since they let us pick, for any index (m,n) displayed by a partner, the index $(i_{m,n},j_{m,n})$ of the best response to that partner in A.

B. Ad-Hoc Adaptation

When faced with an ad-hoc partner, the meta-agent attempts to estimate that partner's niches, potentially across many games, and makes a simple decision at each of its turns:

- If the meta-agent is not confident about its partner's behavioral niche, it chooses actions following generalist policy.
- If the meta-agent is confident that its partner belongs to a certain behavioral niche, it chooses actions following the specialized policy given by the response table to that niche

Algorithm 2 shows a pseudocode of this process, assuming *estimateNiches* returns the estimated niches and some measure of confidence in this estimation, and *getAction* returns the action taken by an agent in a given state.

Algorithm 2: Meta-Agent Action Selection

Input: a game state state, a generalist strategy G, a response table R, a history of interactions with the current partner H, a threshold of confidence C.

Output: the action a to take at the current state $niches, confidence \leftarrow estimateNiches(H);$

 $\begin{array}{ll} \textbf{if} \ confidence > C \ \textbf{then} \\ | \ \ a \leftarrow getAction(G, state) \end{array}$

else

 $specialist \leftarrow R[niches]$ $a \leftarrow getAction(specialist, state)$

 $\mathbf{return}\ a$

In order to this, the meta-agent needs to record information about its partner. In theory, the whole history of moves and game-states could be recorded, but for our purposes we need only to record enough information to reconstruct the metrics we're interested in (e.g. for Communicativeness, the number of hints given and the number of turns with a hint token) and to determine whether we're confident enough to switch from the generalist to the specialized policy.

In our experiment, we use a simple threshold to decide whether to make this switch: if the number of turns played with a certain partner exceeds a fixed threshold, we follow the specialized policy to that partner's apparent niche. Otherwise, we follow the generalist policy.

For our experiments, all information about a partner persists between games. This is accomplished by providing the metaagent with a dummy ID for its current partner at the start of each game. This ID contains no information of strategic importance, but is consistent across games for each partner. This corresponds to the scenario implemented in the Learning Track of the CoG competition [9]. We can also disable this ID if we want the meta-agent to adapt only over the course of a single game (which corresponds to the Mixed track of the competition).

The success of this meta-strategy rest on two assumptions, which we make explicit here:

- That a policy that performs well with agents occupying a given niche in our population will also work well with arbitrary agents (outside that population) that happen to occupy the same niche.
- 2) That it is possible to correctly identify the behavior niche of a partner during gameplay.

Our choice for using average information per play rather than risk aversion as one of our behavioral dimensions can be seen as a step to fulfill assumption 2, since risk-aversion requires estimating the probability that each card is playable, which depends on information our agent will not have access to during game-play (the contents of its own hand).

C. Experimental Set-Up

We implemented three variations of the meta-agent for evaluation purposes:

- An oracle meta-agent, which is given the correct behavioral niche of its partner by a black box module outside the game, and follows the corresponding specialized policy on the response table.
- A generalist meta-agent, which always follows the generalist policy. This is equivalent to setting the threshold to infinite.
- An adaptive meta-agent, which calculates its partner's perceived behavior metrics and always chooses the corresponding specialized response. This corresponds to setting the threshold to zero.

We initially also attempted to set the threshold to other intermediate values, but the agent's performance after reaching the threshold did not vary much for different choices of threshold. For this reason, we restricted ourselves initially to the two extreme values.

We instantiated each of these three meta-agents using the generalist policy and response table from each of the three populations generated by MAP-Elites (for a total of 9 instances). We then evaluated each of these instances by playing 1000 games with each agent for each of the 3 populations (for a total of 27 different experiments).

When evaluating a meta-agent with the same population it was trained with, the oracle and generalist agents serve as a sanity check, as their expected scores can be directly calculated from the match-up data gathered during training. The oracle also acts as an upper bound to the performance of the adaptive agent since, in the best case scenario, the meta-agent will arrive at the same specialized policy as the oracle. The gap in performance between the oracle and the generalist represents the maximum theoretical improvement to be gained by adapting to an appropriate response strategy rather than following a strategy that simply works well on average across all agents. If the adaptive meta-agent manages to achieve better performance than the generalist, it means that this adaptative meta-strategy is paying off in a practical scenario where the niches aren't given.

When evaluating a meta-agent with a population other than the one it was trained with, our objective is to verify to what extent strategies that work for one niche in a certain population also work well for the same niche in a different population. In other words, these experiments verify whether assumption 1 stated in section VI-B holds for this pair of populations.

VII. META-AGENT RESULTS

A. Training Results

To train our meta-agent, we had each agent in the 3 different populations play 400 games with all other agents in the same population. In a population with n covered niches, there are a total of $n^2/2$ possible match-ups, including each agent paired with itself, for a total of $400*n^2/2$ games per population. For each of these match-ups, we record the mean score achieved by that pair of agents. For each agent, we computed their average score across all match-ups, (the intra-population score defined in equation 1). The agent with the highest score constitutes our generalist policy for that population (equation 2). We also computed the best score across all of an agent's match-up, and the indexes of the partner corresponding to that match-up, which makes for the specialized policy for that niche (equation 3).

We call the collection specialized policies for all niches the **response table**. One way to visualize this table is as a collection of line segments on the 2-D plane, mapping the niche (m,n) at one end to its corresponding specialized response $(i_{m,n},j_{m,n})$. Figure 2 shows this visualization for each of the three population. In each graph, every line segment has a blue endpoint, representing a niche, and a red endpoint, representing its specialized response. The niche of the generalist strategy is circled in black. Note that this is not necessarily the niche with the highest number of red endpoints, but the

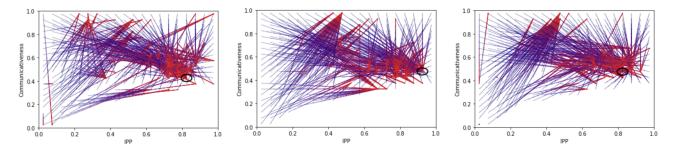


Fig. 2. Visualization of a response table based on populations 1, 2 and 3, from left to right. The blue endpoint of each line segment represents the elite in each niche, and the red endpoint the specialized response to that agent. The circles around (0.825, 0.425), (0.925, 0.475) and (0.825, 0.475) represent the generalist strategy which performs best on average across all pairings in its respective population. Gridlines are omitted so as not to occlude the line segments.

niche with the highest mean score when paired with all the agents in the population.

Note also that the high self-play performance (red) regions of the behavior space shown in figure 1 correspond roughly to the location of the generalist agent and with the regions with higher concentration of red endpoints in figure 2. This suggests that, in general, agents with good self-play performance are also strong when paired with other agents.

However, in all three graphs there are also clusters of red endpoints outside this high-performing region. In particular, all three graphs have a number of red endpoints at or very close to maximum Communicativeness, and their corresponding blue terminations are typically outside the high-performing region. We hypothesize that agents in that high-performing region use mostly similar strategies based on some reasonable assumptions of playability, but agents outside that region, especially agents with IPP ≤ 0.5 tend to play cards "recklessly", whether or not they have information, but still take information into account if it is available. Agents with extreme Communicativeness, which almost never skip a hint, would make ideal partners to these reckless agents.

B. Evaluation Results

Table III shows the results of our evaluations of the metaagent. The average scores for the Oracle, Generalist and Adaptive agent were, respectively, 12.99, 12.25 and 10.37 when evaluated with the same populations they were trained on (underlined entries in the table) and 12.72, 12.04 and 10.35 when evaluated with all agents from all populations.

The gap of around 0.7 point between the Oracle and Generalist scores suggests that there's a potentially significant benefit for correctly choosing a specialized strategy for each partner rather than a generalist strategy. However, the larger gap between the Generalist and Adaptive scores shows that our Adaptive agent falls short of benefiting from this potential adaptation, and one would be better off using the generalist strategy than trying to adapt with our current method.

As we mentioned in section VI-B, the success of the adaptive meta-strategy rested on the assumptions that the choice of the correct adaptive policy is relatively consistent across populations and that we can correctly identify the behavior of a partner during gameplay. For the agents evaluated with

their own respective training populations (underlined entries in table III), the first assumption holds by default, representing an easier scenario for the adaptive agent. Since even on this scenario the Generalist agents significantly outperformed the Adaptive agents, we proceeded to investigate assumption 2, that is, whether the Adaptive agents were correctly identifying its partner's niches.

We investigated this by taking $Adaptive_1$, and having it output at the end of the experiment its estimate for the Communicativeness and IPP scores of its partner, based on the information it stored across all games.

What we found is that, across all 311 partners, the agent's mean absolute error in estimating Communicativeness was 0.11, and the absolute error in estimating IPP was 0.27. The mean *signed* error for these two measurements were -0.07 and -0.20, respectively, meaning that the agent was, on average, underestimating the both behavior metrics of its partner. Since each niche has length 0.05 in each dimension, this means the agent is not assigning the correct niche to its partner most of the time.

We speculate that this is caused by the fact that an agent's behavior isn't specified fully by its own preferences on when to play or hint a card, but also on the actions performed by its partners, which influence the agent's decisions. This would lead to an agent's perceived Communicativeness and IPP to vary drastically between the self-play and ad-hoc play scenarios. Since our Adaptive agent's responses are precomputed based on behavior values measured in self-play, this would lead to not choosing the proper response strategy in the ad-hoc setting.

As an example of how these behavior metrics might differ from self-play to ad-hoc play, consider a hypothetical agent that has some fixed criteria for "likely playability" of a card, and which immediately plays any card deemed likely enough. This agent's IPP in a game would depend on how effective its partner's hints are at establishing playability based on this criteria. However, an arbitrary partner would give different hints than the agent itself would at self-play, so the displayed IPP values would also be different.

	Evaluation population				
Agent	Population 1	Population 2	Population 3		
Oracle ₁	12.84	12.60	12.48		
Oracle ₂	12.60	<u>13.11</u>	12.68		
Oracle ₃	12.49	12.71	13.01		
Generalist ₁	12.12	12.26	12.19		
Generalist ₂	11.75	<u>12.16</u>	11.94		
Generalist ₃	11.68	11.77	12.48		
Adaptive ₁	10.08	9.98	9.85		
Adaptive ₂	10.51	10.47	10.42		
Adaptive ₃	10.67	10.62	10.57		
-	TAB	LE III			

RESULTS OF VALIDATIONS WITH META-AGENT. THE UNDERLINED ENTRIES REPRESENT AGENTS THAT WERE TRAINED AND EVALUATED WITH THE SAME POPULATION

VIII. CONCLUSION AND FUTURE WORK

In this paper, we described our method for generating diverse populations of agents with MAP-Elites using Communicativeness and information per play (IPP) as behavior metrics. We also described a meta-heuristic agent that attempts to adapt during gameplay by using a pre-computed response table indicating which strategy to use when paired with agents in each behavioral niche.

Our results suggest that if the meta-agent is able to correctly identify its partner's niche (as is the case in the Oracle scenarios), it can perform better than if it were to just follow a non-adaptive generalist strategy that maximizes expected performance against a uniform sample of all agents. However, in practice, the meta-agent performs worse than this generalist strategy because it can't correctly identify its partner's niche.

The underlying problem we discovered is to be that an agent's behavior, as captured by our metrics, depends on its partner's behavior. By attempting to measure their behavior during gameplay, we might play in a way as to alter their behavior, interfering with the measurement.

The immediate follow-up question we would like to investigate, then, is how to characterize another agent's behavior during gameplay in a way that avoids this problem. Comparison between the error of estimation of Communicativeness and PIP suggests that some behavior metrics might be more robust to this effect than others, so a possibilty would be to come up with metrics that don't vary as much based on the differences between our own behavior and our partner's self-play behavior. Another possibility would be for the meta-agent to attempt to "triangulate" its partner's behavior by attempting to play various strategies and keeping track of their behavioral response to each strategy.

As another avenue for research, the first assumption discussed in section VI-B rests upon how much variation there is within each niche for agents that are not part of our training population. Some of our results touch on this issue, such as our metrics for diversity between runs in section V-B and the experiments with meta-agents using a different pool for training and evaluation (non-underlined entries in table III). However, more research in necessary, including analysing the response patterns in figure 2 to verify how much variation there is in the best response across runs. It would also be interesting

to compare the pools generated by the current iteration of the MAP-Elites algorithms, using a rule-based representation, to pools generated by other methods of generating agents, such as reinforcement learning and neuroevolution [29].

Finally, it would be interesting to verify whether our characterization of Hanabi players based on the behavior metrics described in this paper is useful to describe the behavior of human players. Some questions we would like to address are: what niches are predominantly occupied by humans? Do our current behavior metrics capture the variations of play exhibited by humans? Is there a behavioral region that plays best with humans? Do the answers to these questions depend on a particular human's skill level or some other factor? Data could be gathered from new experiments addressing these questions or from from previous studies involving human players such as [7], [21], [23].

ACKNOWLEDGMENT

Rodrigo Canaan, Andy Nealen and Julian Togelius gratefully acknowledge the financial support from Honda Research Institute Europe (HRI-EU).

REFERENCES

- [1] P. Stone, G. A. Kaminka, S. Kraus, and J. S. Rosenschein, "Ad hoc autonomous agent teams: Collaboration without pre-coordination," in *Twenty-Fourth AAAI Conference on Artificial Intelligence*, 2010.
- [2] D. Premack and G. Woodruff, "Does the chimpanzee have a theory of mind?" *Behavioral and brain sciences*, vol. 1, no. 4, pp. 515–526, 1978.
- [3] J.-B. Mouret and J. Clune, "Illuminating search spaces by mapping elites," arXiv preprint arXiv:1504.04909, 2015.
- [4] R. Canaan, J. Togelius, A. Nealen, and S. Menzel, "Diverse agents for ad-hoc cooperation in hanabi," in 2019 IEEE Conference on Games (CoG). IEEE, 2019, pp. 1–8.
- [5] H. Osawa, "Solving hanabi: Estimating hands by opponent's actions in cooperative game with incomplete information." in AAAI workshop: Computer Poker and Imperfect Information, 2015, pp. 37–43.
- [6] M. J. van den Bergh, A. Hommelberg, W. A. Kosters, and F. M. Spieksma, "Aspects of the cooperative card game hanabi," in *Benelux Conference on Artificial Intelligence*. Springer, 2016, pp. 93–105.
- [7] M. Eger, C. Martens, and M. A. Córdoba, "An intentional ai for hanabi," in *Computational Intelligence and Games (CIG)*, 2017 IEEE Conference on. IEEE, 2017, pp. 68–75.
- [8] J. Walton-Rivers, P. R. Williams, R. Bartle, D. Perez-Liebana, and S. M. Lucas, "Evaluating and modelling hanabi-playing agents," in Evolutionary Computation (CEC), 2017 IEEE Congress on. IEEE, 2017, pp. 1382–1389.
- [9] J. Walton-Rivers, P. R. Williams, and R. Bartle, "The 2018 hanabi competition," in 2019 IEEE Conference on Games (CoG). IEEE, 2019, pp. 1–8.
- [10] C. B. Browne, E. Powley, D. Whitehouse, S. M. Lucas, P. I. Cowling, P. Rohlfshagen, S. Tavener, D. Perez, S. Samothrakis, and S. Colton, "A survey of monte carlo tree search methods," *IEEE Transactions on Computational Intelligence and AI in games*, vol. 4, no. 1, pp. 1–43, 2012.
- [11] J. Goodman, "Re-determinizing information set monte carlo tree search in hanabi," arXiv preprint arXiv:1902.06075, 2019.
- [12] J. Walton-Rivers, "Hanabi competition results," https://community. fossgalaxy.com/t/hanabi-competition-results/154, 2018, access: 05/14/2018.
- [13] R. Canaan, H. Shen, R. Torrado, J. Togelius, A. Nealen, and S. Menzel, "Evolving agents for the hanabi 2018 cig competition," in 2018 IEEE Conference on Computational Intelligence and Games (CIG). IEEE, 2018, pp. 1–8.
- [14] N. Bard, J. N. Foerster, S. Chandar, N. Burch, M. Lanctot, H. F. Song, E. Parisotto, V. Dumoulin, S. Moitra, E. Hughes *et al.*, "The hanabi challenge: A new frontier for ai research," *arXiv preprint arXiv:1902.00506*, 2019.

- [15] M. Hessel, J. Modayil, H. Van Hasselt, T. Schaul, G. Ostrovski, W. Dabney, D. Horgan, B. Piot, M. Azar, and D. Silver, "Rainbow: Combining improvements in deep reinforcement learning," in *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [16] A. Lerer, H. Hu, J. Foerster, and N. Brown, "Improving policies via search in cooperative partially observable games," arXiv preprint arXiv:1912.02318 to be presented at AAAI 2020, 2019.
- [17] H. Hu and J. N. Foerster, "Simplified action decoder for deep multiagent reinforcement learning," arXiv preprint arXiv:1912.02288 to be presented at ICLR 2020, 2019.
- [18] C. Cox, J. De Silva, P. Deorsey, F. H. Kenter, T. Retter, and J. Tobin, "How to make the perfect fireworks display: Two strategies for hanabi," *Mathematics Magazine*, vol. 88, no. 5, pp. 323–336, 2015.
- [19] B. Bouzy, "Playing hanabi near-optimally," in Advances in Computer Games. Springer, 2017, pp. 51–62.
- [20] J. Wu, "State of the art hanabi bots + simulation framework in rust," https://github.com/WuTheFWasThat/hanabi.rs, 2016, access: 05/14/2018.
- [21] C. Liang, J. Proft, E. Andersen, and R. A. Knepper, "Implicit communication of actionable information in human-ai teams," in *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, 2019, pp. 1–13.
- [22] H. P. Grice, "Logic and conversation," in Speech acts. Brill, 1975, pp. 41–58.
- [23] H. Hu, A. Lerer, A. Peysakhovich, and J. Foerster, "" other-play" for zero-shot coordination," *arXiv preprint arXiv:2003.02979*, 2020.
- [24] J. K. Pugh, L. B. Soros, and K. O. Stanley, "Quality diversity: A new frontier for evolutionary computation," *Frontiers in Robotics and AI*, vol. 3, p. 40, 2016.
- [25] J. Lehman and K. O. Stanley, "Abandoning objectives: Evolution through the search for novelty alone," *Evolutionary computation*, vol. 19, no. 2, pp. 189–223, 2011.
- [26] K. Deb, "Multi-objective optimization," in Search methodologies. Springer, 2014, pp. 403–449.
- [27] A. Cully, J. Clune, D. Tarapore, and J.-B. Mouret, "Robots that can adapt like animals," *Nature*, vol. 521, no. 7553, p. 503, 2015.
- [28] R. W. Hamming, "Error detecting and error correcting codes," The Bell system technical journal, vol. 29, no. 2, pp. 147–160, 1950.
- [29] F. P. Such, V. Madhavan, E. Conti, J. Lehman, K. O. Stanley, and J. Clune, "Deep neuroevolution: Genetic algorithms are a competitive alternative for training deep neural networks for reinforcement learning," arXiv preprint arXiv:1712.06567, 2017.