# On Emergent Communication in Competitive Multi-Agent Teams

Paul Pu Liang, Jeffrey Chen, Ruslan Salakhutdinov, Louis-Philippe Morency, Satwik Kottur

Carnegie Mellon University

pliang@cs.cmu.edu

## ABSTRACT

Several recent works have found the emergence of grounded compositional language in the communication protocols developed by mostly cooperative multi-agent systems when learned end-to-end to maximize performance on a downstream task. However, human populations learn to solve complex tasks involving communicative behaviors not only in fully cooperative settings but also in scenarios where competition acts as an additional external pressure for improvement. In this work, we investigate whether competition for performance from an external, similar agent team could act as a social influence that encourages multi-agent populations to develop better communication protocols for improved performance, compositionality, and convergence speed. We start from *Task & Talk*, a previously proposed referential game between two cooperative agents as our testbed and extend it into *Task, Talk & Compete*, a game involving two competitive teams each consisting of two aforementioned cooperative agents. Using this new setting, we provide an empirical study demonstrating the impact of competitive influence on multi-agent teams. Our results show that an external competitive influence leads to improved accuracy and generalization, as well as faster emergence of communicative languages that are more informative and compositional.

## CCS CONCEPTS

• **Theory of computation** → *Multi-agent learning*; • **Computing methodologies** → *Artificial intelligence*; *Multi-agent systems*;

## KEYWORDS

Learning agent-to-agent interactions; Multi-agent learning

## 1 INTRODUCTION

The emergence and evolution of languages through human life, societies, and cultures has always been one of the hallmarks of human intelligence [11, 34, 42]. Humans intelligently communicate through language to solve multiple real-world tasks involving vision, navigation, reasoning, and learning [35, 43]. Language emerges naturally for us humans in both individuals through inner speech [1, 3, 5] as well as in groups through grounded dialog in
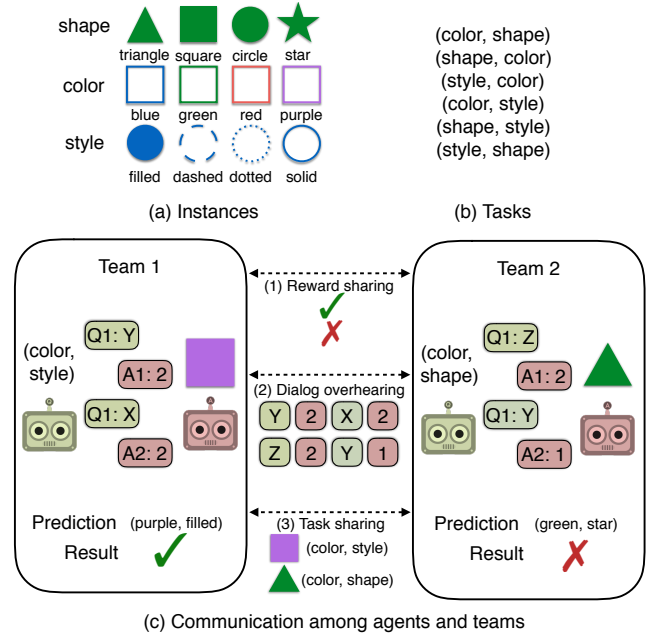
Figure 1: We propose the *Task, Talk & Compete* game involving two competitive teams each consisting of multiple rounds of dialog between Q-BOT and A-BOT. Within each team, A-BOT is given a target instance unknown to Q-BOT and Q-BOT is assigned a task (e.g. find the color and shape of the instance), unknown to A-BOT, to uncover certain attributes of the target instance. This informational asymmetry necessitates communication between the two agents via multiple rounds of dialog where Q-BOT asks questions regarding the task and A-BOT provides answers using the target instance. We investigate whether *competition for performance* from an external, similar agent team 2 could result in improved compositionality of emergent language and convergence of task performance within team 1. Competition is introduced through three aspects: 1) *reward sharing*, 2) *dialogue overhearing*, and 3) *task sharing*. Our hypothesis is that teams are able to leverage information from the performance of the other team and learn more efficiently beyond the sole reward signal obtained from their own performance. Our findings show that competition for performance with a similar team leads to improved overall accuracy of generalization as well as faster emergence of emergent language.

both cooperative and competitive settings [23, 24, 41]. As a result, there has been a push to build artificial intelligence models that can effectively communicate their own intentions to both humans as well as other agents through language and dialog grounded in vision [10, 29, 30, 38] and navigation [14].

A recent line of work has applied reinforcement learning techniques [44] for end-to-end learning of communication protocols between agents situated in virtual environments [9, 27] and found the emergence of grounded, interpretable, and compositional symbolic language in the communication protocols developed by the agents [2, 26]. To bring this closer to how natural language evolves within communities, we observe that human populations learn to solve complex tasks involving communicative behaviors not only in fully cooperative settings but also in scenarios where competition acts as an additional external pressure for improvement [15, 33]. Furthermore, recent studies have also provided support for external competition in behavioral economics [17] and evolutionary biology [12, 28, 36]. Inspired by the emergence of human behavior and language in competitive settings, the goal of our paper is to therefore understand how social influences like competition affects emergent language in multi-agent communication.

In this work, we conceptually distinguish two types of competition: 1) constant sum competition, where agents are competing for a finite amount of resources and thus a gain for a group of agents results in a loss for another [4, 32], and 2) variable sum or non-constant sum competition, where competition is inherent but both mutual gains and mutual losses of power are possible [6, 40]. For the sake of simplicity, we call these *competition for resources* and *competition for performance* respectively and focus on the latter. In this setting, each individual (or team) performs similar tasks with access to individual resources but is motivated to do better due to the external pressure from seeing how well other individuals (or teams) are performing. We hypothesize that *competition for performance* from an external, similar agent team could act as a social influence that encourages multi-agent populations to develop better communication protocols for improved performance, compositionality, and convergence.

To investigate this hypothesis, we extend the *Task & Talk* referential game between two cooperative agents [25] into *Task, Talk & Compete*, a game involving two competitive teams each consisting of the aforementioned two cooperative agents (Figure 1). At a high level, *Task & Talk* requires the two cooperative agents (Q-BOT and A-BOT) to solve a *task* by interacting with each other via *dialog*, at the end of which a *reward* is assigned to measure their performance. In such a setting, we introduce competition for performance between two teams of Q-BOT and A-BOT through three aspects: 1) *reward sharing*, where we modify the reward structure so that a team can assess their performance relative to the other team, 2) *dialog overhearing*, where we modify the agent's policy networks such that they can overhear the concurrent dialog from the other team, and 3) *task sharing*, where teams gain additional information about the tasks given to the opposing team. Using this new setting, we provide an empirical study demonstrating the impact of competitive influence on multi-agent dialog games. Our results show that external competitive influence leads to an improved accuracy and generalization, as well as faster emergence of communicative languages that are more informative and compositional.

## 2 THE *TASK & TALK* GAME

The *Task & Talk* benchmark is a cooperative reference game proposed to evaluate the emergence of language in multi-agent populations with imperfect information [13, 25]. *Task & Talk* takes place between two agents Q-BOT (which is tasked to ask questions) and A-BOT (which is tasked to answer questions) in a synthetic world of instances comprised of three attributes: color, style, and shape. At the start of the game, A-BOT is given a target instance $I$ unknown to Q-BOT and Q-BOT is assigned a task $G$ (unknown to A-BOT) to uncover certain attributes of the target instance $I$. This informational asymmetry necessitates communication between the two agents via multiple rounds of dialog where Q-BOT asks questions regarding the task and A-BOT provides answers using the target instance. At the end of the game, Q-BOT uses the information conveyed from the dialog with A-BOT to solve the task at hand. Both agents are given equal rewards based on the accuracy of Q-BOT's prediction.[1]

**Base States and Actions:** Each agent begins by observing its specific input: task $G$ for Q-BOT and instance $I$ for A-BOT. The game proceeds as a dialog over multiple rounds $t = 1, ..., T$. We use lower case characters (e.g. $s_t^Q$) to denote token symbols and upper case $S_t^Q$ to denote corresponding representations. Q-BOT is modeled with three modules – speaking, listening, and prediction. At round $t$, Q-BOT stores an initial state representation $S_{t-1}^Q$ from which it conditionally generates output utterances $q_t \in V_Q$ where $V_Q$ is the vocabulary of Q-BOT. $S_{t-1}^Q$ is updated using answers $a_t$ from A-BOT and is used to make a prediction $\hat{w}_G$ in the final round. A-BOT is modeled with two modules – speaking and listening. A-BOT encodes instance $I$ into its initial state $S_t^A$ from which it conditionally generates output utterances $a_t \in V_A$ where $V_A$ is the vocabulary of Q-BOT. $S_t^A$ is updated using questions $q_t$ from Q-BOT.

In more detail, Q-BOT and A-BOT are modeled as **stochastic policies** $\pi_Q(q_t|s_t^Q; \theta_Q)$ and $\pi_A(a_t|s_t^A; \theta_A)$ implemented as recurrent networks [20, 21]. At the beginning of round $t$, Q-BOT observes state $s_t^Q = [G, q_1, a_1, \ldots, q_{t-1}, a_{t-1}]$ representing the task $G$ and the dialog conveyed up to round $t-1$. Q-BOT conditions on state $s_t^Q$ and utters a question represented by some token $q_t \in V_Q$. A-BOT also observes the dialog history and this new utterance as state $s_t^A = [I, q_1, a_1, \ldots, q_{t-1}, a_{t-1}, q_t]$. A-BOT conditions on this state $s_t^A$ and utters an answer $a_t \in V_A$. At the final round, Q-BOT predicts a pair of attributes $\hat{w}^G = (\hat{w}_1^G, \hat{w}_2^G)$ using a network $\pi_G(g|s_T^Q; \theta_Q)$ to solve the task.

**Learning the Policy**: Q-BOT and A-BOT are trained to cooperate by maximizing a shared objective function that is determined by whether Q-BOT is able solve the task at hand. Q-BOT and A-BOT receive an identical **base reward** of $R$ if Q-BOT's prediction $\hat{w}^G$ matches ground truth $w^G$ and a negative reward of $-10R$ otherwise. $R$ is a hyperparameter that affects the rate of convergence.

$$J(\theta_Q, \theta_A) = \mathbb{E}_{\pi_Q, \pi_A}\left[\mathcal{R}(\hat{w}^G, w^G)\right] \quad (1)$$

where $\mathcal{R}$ is a reward function. To train these agents, we learn policy parameters $\theta_Q$ and $\theta_A$ that maximize the expected reward $J(\theta_Q, \theta_A)$

---

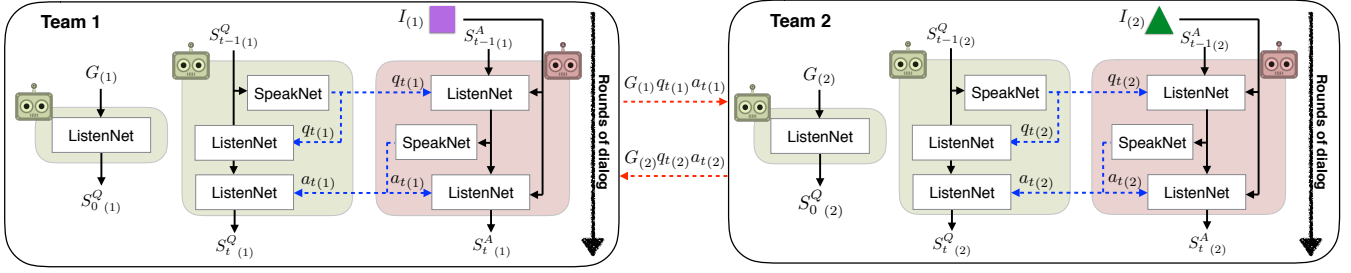[1] detailed review of *Task & Talk* in the appendix.

Figure 2: Detailed policy networks (implemented as conditional recurrent networks) and communication protocols for Q-BOT and A-BOT in both teams. At the start, A-BOT is given a target instance $I$ unknown to Q-BOT and Q-BOT is assigned a task $G$ (unknown to A-BOT) to uncover certain attributes of the target instance $I$. At each round $t$, Q-BOT observes state $s_t^Q$ and utters some token $q_t \in V_Q$. A-BOT observes the history and this new utterance as state $s_t^A$ and utters $a_t \in V_A$. At the final round, Q-BOT predicts a pair of attribute values $\hat{w}^G = (\hat{w}_1^G, \hat{w}_2^G)$ to solve the task. Policy parameters $\theta_Q$ and $\theta_A$ are updated using the REINFORCE policy gradient algorithm [44] on a positive reward $+R$ if Q-BOT's prediction $\hat{w}^G$ matches ground truth $w^G$ and a negative reward of $-10R$ otherwise. Competition is introduced via 1) *reward sharing* ($R_{(2)}$ passed to team 1), 2) *dialog overhearing* (team 1 overhears and conditions on $q_{t(2)}, a_{t(2)}$), and 3) *task sharing* (team 1 knows about $I_{(2)}$ and $G_{(2)}$). Blue dotted lines represent (cooperative) dialog within a team and red dotted lines represent (competitive) dialog across teams.

---

**Algorithm 1** Training a single team of cooperative agents.

**COOPERATIVETRAIN:**

1: Given Q-BOT params $\theta_Q$ and A-BOT params $\theta_A$.
2: **for** $(I, G)$ in each batch **do**
3:   **for** communication round $t = 1, \ldots, T$ **do**
4:     $s_t^Q = [G, q_1, a_1, \ldots, q_{t-1}, a_{t-1}]$
5:     $q_t = \pi_Q(q_t | s_t^Q; \theta_Q)$
6:     $s_t^A = [I, q_1, a_1, \ldots, q_{t-1}, a_{t-1}, q_t]$
7:     $a_t = \pi_A(a_t | s_t^A; \theta_A)$
8:   **end for**
9:   $\hat{w}^G = (\hat{w}_1^G, \hat{w}_2^G) = \pi_G(g | s_T^Q; \theta_Q)$
10:   $J(\theta_Q, \theta_A) = \mathbb{E}_{\pi_Q, \pi_A} \left[ \mathcal{R}(\hat{w}^G, w^G) \right]$.
11:   $\theta_Q = \theta_Q + \eta \nabla_{\theta_Q} J(\theta_Q, \theta_A)$.
12:   $\theta_A = \theta_A + \eta \nabla_{\theta_A} J(\theta_Q, \theta_A)$.
         ▷ Update $\theta_Q, \theta_A$ using estimated $\nabla J(\theta_Q, \theta_A)$.
13: **end for**

---

using the REINFORCE policy gradient algorithm [44]. The expectation of policy gradients for $\theta_Q$ and $\theta_A$ are given by

$$\nabla_{\theta_Q} J(\theta_Q, \theta_A) = \mathbb{E}_{\pi_Q, \pi_A} \left[ \mathcal{R}(\hat{w}^G, w^G) \nabla_{\theta_Q} \log \pi_Q \left( q_t | s_t^Q; \theta_Q \right) \right], \quad (2)$$

$$\nabla_{\theta_A} J(\theta_Q, \theta_A) = \mathbb{E}_{\pi_Q, \pi_A} \left[ \mathcal{R}(\hat{w}^G, w^G) \nabla_{\theta_A} \log \pi_A \left( a_t | s_t^A; \theta_A \right) \right]. \quad (3)$$

These expectation are approximated by sample averages across object instances and tasks in a batch, as well as across dialog rounds for a given object instance and task. Using the estimated gradients, the parameters $\theta_Q$ and $\theta_A$ are updated using gradient-based methods in an alternating fashion until convergence. We summarize the procedure for cooperative training for a single team of agents in Algorithm 1, which we call COOPERATIVETRAIN.

## 3  THE *TASK, TALK & COMPETE* GAME

We modify the *Task & Talk* benchmark into the *Task, Talk & Compete* game (Figure 2). Our setting now consists of two teams of agents: Q-BOT$_{(1)}$ and A-BOT$_{(1)}$ belonging to team 1 and Q-BOT$_{(2)}$ and A-BOT$_{(2)}$ belonging to team 2. Similar to the traditional *Task & Talk* game, the agents Q-BOT$_{(1)}$ and A-BOT$_{(1)}$ in team 1 cooperate and communicate to solve a task, and likewise for the agents in team 2.

The *Task, Talk & Compete* game begins with two target instances $I_{(1)}$ and $I_{(2)}$ presented to Q-BOT$_{(1)}$ and Q-BOT$_{(2)}$ respectively, and two tasks $G_{(1)}$ and $G_{(2)}$ presented to A-BOT$_{(1)}$ and A-BOT$_{(2)}$ respectively. Within a team, we largely follow the setting in *Task & Talk* as described in the previous section. A *team* consists of a pair agents Q-BOT and A-BOT cooperating in a partially observable world to solve task $G$ given instance $I$. The key difference is that agents in one team are *competing* against those in the other team. In the following subsections, we explain the various sources of competition that we introduce into the game and the modified training procedure for teams of agents.

We use subscripts to index the rounds and subscripts in parenthesis to index which team the agents belong to (i.e. $s_{t\ (1)}^Q$). We drop the team subscript if it is clear from the context (i.e. same team). Note that for grounding to happen across teams (i.e. team 1 to understand team 2 during dialog overhearing and vice-versa), the teams *share vocabularies sizes* (i.e. $V_Q$ is shared by Q-BOT$_{(1)}$ and Q-BOT$_{(2)}$). They do not share vocabularies since both teams train differently but they have access to the same number of symbols.

### 3.1  Sources of Competition

When the two teams do not share any information and are trained completely independently, the *Task, Talk & Compete* game reduces to (two copies of) the *Task & Talk* game. Therefore, information sharing across the teams is necessary to introduce competition. In the following section we highlight the information sharing that can happen in *Task, Talk & Compete* and the various sources of competition that can subsequently arise.

**Algorithm 2** Training two teams to compete against each other.

---

**COMPETITIVETRAIN:**

1: Given: Q-BOT$_{(j)}$ params $\theta_{Q(j)}$; A-BOT$_{(j)}$ params $\theta_{A(j)}$, for $j = 0, 1$.
2: **for** $(I, G)$ in each batch **do**
3:    **for** communication round $t = 1, \ldots, T$ **do**
4:       **for** $j = 0, 1$ **do**
5:          $j' = 1 - j$     ▷ Index of the other team
6:          $s^Q_{t\,(j)} = [G_{(j)}, q_{1(j)}, a_{1(j)}, \ldots, q_{t-1(j)}, a_{t-1(j)}]$
7:          **if** <u>Dialog Overhearing</u> **then**
8:             $s^Q_{t\,(j)}+ = [q_{1(j')}, a_{1(j')}, \ldots, q_{t-1(j')}, a_{t-1(j')}]$
                ▷ Overhear dialog from the other team
9:          **end if**
10:         **if** <u>Task Sharing</u> **then**
11:            $s^Q_{t\,(j)}+ = [G_{(j')}]$
12:         **end if**
13:         $q_{t(j)} = \pi_Q(q_t | s^Q_{t\,(j)}; \theta_{Q(j)})$
14:         (and symmetrically for $s^A_{t\,(j)}$ of A-BOT$_{(j)}$)
15:       **end for**
16:    **end for**
17:    Compute $\hat{w}^G = (\hat{w}^G_1, \hat{w}^G_2) = \pi_G(g | s^Q_T; \theta_Q)$ independently for each team.
18:    **if** <u>Reward Sharing</u> **then**
19:       $J(\theta_{Q(j)}, \theta_{A(j)}) = \mathbb{E}_{\pi_Q, \pi_A}\left[\mathcal{R}_{\text{shared}}(\hat{w}^G_{(j)}, w^G_{(j)})\right], j = 0, 1.$
         ▷ Compute $J()$ jointly using predictions of *both* teams
20:    **else**
21:       $J(\theta_{Q(j)}, \theta_{A(j)}) = \mathbb{E}_{\pi_Q, \pi_A}\left[\mathcal{R}(\hat{w}^G_{(j)}, w^G_{(j)})\right], j = 0, 1.$
         ▷ Compute $J()$ independently for *each* team
22:    **end if**
23:    $\theta_{Q(j)} = \theta_{Q(j)} + \eta \nabla_{\theta_{Q(j)}} J(\theta_{Q(j)}, \theta_{A(j)}), j = 0, 1.$
24:    $\theta_{A(j)} = \theta_{A(j)} + \eta \nabla_{\theta_{A(j)}} J(\theta_{Q(j)}, \theta_{A(j)}), j = 0, 1.$
         ▷ Update parameters using estimated $\nabla J()$.
25: **end for**

**FULLTRAIN:**

26: **for** $epoch = 1, \ldots, \text{MAX}$ **do**
27:    COOPERATIVETRAIN agents within team 1.
28:    COOPERATIVETRAIN agents within team 2.
29:    COMPETITIVETRAIN agents across teams 1 and 2.
30: **end for**

---

**Reward Sharing (RS):** We modify the reward structure so that a team can assess their performance relative to other teams. Starting with a base reward of $R$ in the fully cooperative setting, we modify this reward into the competitive setting in the following Table:

|  | Team 2 ✓ | Team 2 ✗ |
|---|---|---|
| Team 1 ✓ | $(+R, +R)$ | $(+R, -100R)$ |
| Team 1 ✗ | $(-100R, +R)$ | $(-10R, -10R)$ |

When both teams get the same result this reduces to the base reward setting: $+R$ for correct answers and $-10R$ for wrong ones. When there is asymmetry in performance across the two teams, the correct team gains a reward of $+R$ while the incorrect team suffers a larger penalty of $-100R$. This encourages teams to compete and assess their performance relative to the other team.

**Dialog Overhearing (DO):** Overhearing the conversations of another team is a common way to get secret information about the tactics, knowledge, and progress of one's competitors. In a similar fashion, we modify the agent's policy networks such that they can now overhear the concurrent dialog from other teams. Take Q-BOT$_{(1)}$ and A-BOT$_{(1)}$ in team 1 for example. At round $t$, Q-BOT$_{(1)}$ now observes state

$$s^Q_{t\,(1)} = [G_{(1)}, q_{1(1)}, a_{1(1)}, q_{1(2)}, a_{1(2)}, \ldots,$$
$$q_{t-1(1)}, a_{t-1(1)}, q_{t-1(2)}, a_{t-1(2)}], \tag{4}$$

and similarly for Q-BOT$_{(2)}$. A-BOT$_{(1)}$ observes

$$s^A_{t\,(1)} = [I_{(1)}, q_{1(1)}, a_{1(1)}, q_{1(2)}, a_{1(2)}, \ldots,$$
$$q_{t-1(1)}, a_{t-1(1)}, q_{t-1(2)}, a_{t-1(2)}, q_{t(1)}, q_{t(2)}]. \tag{5}$$

and similarly for A-BOT$_{(2)}$. We view this as augmenting reward sharing by informing the agents as to why they were penalized: a team can listen to what the other team is communicating and use that to its own advantage. In practice, we define an overhear fraction $\rho$ which determines how often overhearing occurs during the training epochs. $\rho$ is a hyperparameter in our experiments.

**Task Sharing (TS):** Finally, we fully augment both reward sharing and dialog overhearing with task sharing, where the two sets of instances and tasks are known to both teams (i.e. $G_{(2)}$ is added to Equation 4 and $I_{(2)}$ is added to Equation 5). Task sharing adds an additional level of grounding: one team can now ground the overheard dialog in a specific task and compare their performance to that of the other team. Overall, Algorithm 2 summarizes the procedure for training teams of agents in a competitive setting using the aforementioned three sources of competition. We call the resulting algorithm COMPETITIVETRAIN.

**Asynchronous Training:** The baseline from Kottur et al. [25] only evaluates the performance of a single team at test time. For fair comparison with this baseline, we evaluate performance with a single team 1 as well. In the case of dialog overhearing and task sharing, we replace the overheard symbols and tasks from task 2 with zeros during evaluation. This removes the confounding explanation that improved performance is due to having more information from team 2 during testing.

To prevent data mismatch during training and testing [16, 39], we use a three-stage training process. During stage (1), Q-BOT$_{(1)}$ and A-BOT$_{(1)}$ are trained to cooperate, independent of team 2. All information shared from team 2 is passed to team 1 as zeros. During stage (2), Q-BOT$_{(2)}$ and A-BOT$_{(2)}$ are trained to cooperate, and in stage (3), both teams are trained together with reward sharing and/or dialog overhearing and/or task sharing. These three stages are repeated until convergence. Intuitively, this procedure means that the agents within each team must learn how to cooperate with each other in addition to competing with the other team. The final algorithm which takes into account asynchronous training both within and across teams is shown in Algorithm 2, which we call FULLTRAIN. FullTrain is the final algorithm used for training both teams of agents.

Due to how we model competition in our setting (either through reward sharing or task and dialog overhearing), our hypothesis is that teams are able to leverage information from the performance of the other team and learn more efficiently beyond the signal obtained solely from their performance.

## 4 EXPERIMENTAL SETUP

We present our experimental results and observations on the effect of competitive influence on both final task performance and emergence of grounded, interpretable, and compositional language. Our experimental testbed is the *Task, Talk & Compete* game. We implement a variety of different algorithms spanning cooperative baselines and competitive methods which we will detail in the following subsections.[2]

### 4.1 Cooperative Baselines

These are baseline methods that involve only cooperative team (or teams) of agents. These baselines test for various confounders of our experiments (e.g. improved performance due to increase in the number of parameters).

(1) COOP, BASE: a single team, fully cooperative setting with reward structure $(+R, -10R)$, which is the baseline from [25].

(2) COOP, REWARDS: a single team, fully cooperative setting with reward structure $(+R, -100R)$ adjusted for our "extreme" reward setting that is more strict in penalizing incorrect answers. This baseline ensures that our improvements are not simply due to better reward shaping [18, 19].

(3) COOP, PARAMS: a single team, fully cooperative setting with (roughly) double the number of parameters as compared to the baseline. This ensures that the improvement in performance is not due to an increase in the number of parameters during dialog overhearing or task sharing. This baseline is obtained by increasing the LSTM hidden size for question and answer bots such that the total number of parameters match the competitive settings with roughly double the number of parameters.

(4) COOP, DOUBLE: two teams each trained independently without sharing any information. At test time, the team with the higher validation score is used to evaluate performance. This baseline ensures that our improvements are not simply due to double the chance of beginning with a better random seed or luckier training.

### 4.2 Competitive Methods

The following methods introduce an extra level of competition between teams on top of cooperation within teams. Competitive behavior is encouraged via combinations of reward sharing, dialog overhearing, and task sharing.

(1) COMP, TS: two competitive teams with task sharing.

(2) COMP, DO: two competitive teams with dialog overhearing.

(3) COMP, DO+TS: two competitive teams with dialog overhearing and task sharing.

(4) COMP, RS: two competitive teams with reward sharing.

(5) COMP, RS+TS: two competitive teams with reward sharing and task sharing.

(6) COMP, RS+DO: two competitive teams with reward sharing and dialog overhearing.

(7) COMP, RS+DO+TS: two competitive teams with reward sharing, dialog overhearing, and task sharing.

**Hyperparameters:** We perform all experiments with the same hyperparameters (except LSTM hidden dimension which is fixed at 100 but increased to 150 for the COOP, PARAMS setting to experiment with an increase in the number of parameters). We set the reward multiplier $R = 100$, overhear fraction $\rho = 0.5$, and vocabulary sizes of Q-BOT and A-BOT to be $|V_Q| = 3$ and $|V_A| = 4$ respectively. Following Kottur et al. [25], we set A-BOT$_{(1)}$ and A-BOT$_{(2)}$ to be memoryless to ensure consistent A-BOT grounding across rounds which is important for generalization and compositional language. All other parameters follow those in Kottur et al. [25].

**Metrics:** In addition to evaluating the train and test accuracies [25], we would also like to investigate the impact of competitive pressure on the emergence of language among agents. We measure **Instantaneous Coordination (IC)** [22] defined as the mutual information between one agent's message and *another* agent's action, i.e. MI$(a_{t(j)}, \hat{w}^G_{i(j)})$. Higher IC implies that Q-BOT's action depends more strongly on A-BOT's message (and vice versa), which is in turn indicative that messages are used in positive manner to signal actions within a team. Although Lowe et al. [31] mentioned other metrics such as speaker consistency [8, 22] and entropy, we believe that IC is the most suited for question answering dialog tasks where responding to *another* agent's messages is key. We also measured several other metrics that were recently proposed to measure how informative a language is with respect to the agent's actions [31]: **Speaker Consistency (SC)** measures the mutual information between an agent's message and its own future action: MI$(q_{t(j)}, \hat{w}^G_{i(j)})$ and **Entropy (H)** which measures the entropy of an agent's sequence of outgoing messages.

**Training Details:** For cooperative baselines, we set the maximum number of epochs to be 100,000 and stop training early when training accuracy reaches 100%. For competitive baselines, we also set the maximum number of epochs to be 100,000 and stop training early when *the first team* reaches a training accuracy of 100%. This team is labeled as the *winning* team and is the focus of our experiments. The other *losing* team can be viewed as an auxiliary team that helps the performance of the winning team. All experiments are repeated 10 times with randomly chosen random seeds. Results are reported as average ± standard deviation over the 10 runs. Implementation details are provided in the appendix.

## 5 RESULTS AND ANALYSIS

We study the effect of competition between teams on 1) the generalization abilities of the agents in new environments during test-time, 2) the rate of convergence of train and test accuracies during training, and 3) the emergence of informative communication protocols between agents when solving the *Task, Talk & Compete* game.

### 5.1 Qualitative Results

We begin by studying the effect of competition on the performance of agents in the *Task, Talk & Compete* game. The teams are trained in various cooperative and competitive settings. We aggregate scores for both winning and losing teams as described in the training details above. These results are summarized in Table 1. From our results, we draw the following general observations regarding the generalization capabilities of the trained agents:

---

| Type | Method | RS | DO | TS | Train Acc (%) | | Test Acc (%) | |
|---|---|---|---|---|---|---|---|---|
| | | | | | Winning Team | Losing Team | Winning Team | Losing Team |
| Cooperative baselines | Coop, base [25] | ✗ | ✗ | ✗ | 88.5 ± 11.6 | - | 45.6 ± 18.9 | - |
| | Coop, rewards [18, 19] | ✗ | ✗ | ✗ | 87.0 ± 13.7 | - | 49.7 ± 22.9 | - |
| | Coop, params | ✗ | ✗ | ✗ | 85.5 ± 14.6 | - | 53.3 ± 26.2 | - |
| | Coop, double | ✗ | ✗ | ✗ | 91.4 ± 12.0 | 74.6 ± 16.6 | 57.8 ± 28.5 | 38.3 ± 27.0 |
| Competitive methods | Comp, TS | ✗ | ✗ | ✓ | 94.4 ± 3.0 | 80.4 ± 13.3 | 53.1 ± 19.5 | 41.7 ± 25.7 |
| | Comp, DO | ✗ | ✓ | ✗ | 96.8 ± 5.4 | 71.4 ± 14.2 | 65.7 ± 26.1 | 23.1 ± 10.7 |
| | Comp, DO+TS | ✗ | ✓ | ✓ | 98.6 ± 2.0 | 75.8 ± 12.1 | **75.8 ± 17.9** | 41.1 ± 26.0 |
| | Comp, RS | ✓ | ✗ | ✗ | 99.5 ± 1.3 | 79.6 ± 12.2 | 68.5 ± 19.4 | 43.8 ± 20.9 |
| | Comp, RS+TS | ✓ | ✗ | ✓ | 99.5 ± 1.4 | 66.6 ± 10.3 | 68.9 ± 16.9 | 26.7 ± 11.3 |
| | Comp, RS+DO | ✓ | ✓ | ✗ | **100.0 ± 0.0** | 63.9 ± 12.8 | **78.3 ± 10.7** | 28.3 ± 14.4 |
| | Comp, RS+DO+TS | ✓ | ✓ | ✓ | 98.8 ± 2.4 | 78.9 ± 10.9 | **77.2 ± 16.5** | 28.9 ± 21.8 |

Table 1: Train and test accuracies measured across teams trained in various cooperative and competitive settings. All cooperative baselines are in shades of blue and competitive teams are in red. RS: reward sharing, DO: dialog overhearing, TS: task sharing. Accuracies are reported separately for winning and losing teams with best accuracies for winning teams in bold. Winning teams in competitive settings display faster convergence and improved performance.

(1) **Sharing messages via overhearing dialog improves generalization performance:** Dialog overhearing contributes the most towards improvement in test accuracy, from below 60% in the baselines without to 75.8% with dialog overhearing. We believe this is because dialog overhearing transmits the most amount of information to the other team, as compared to a single scalar in reward sharing or a single image in task sharing.

(2) **Composing sources of competition improves performance:** While dialog overhearing on its own displays strong improvements in test performance, we found that composing multiple sources of competition improves performance even more. In the Comp, RS+DO setting, the winning team's train accuracy quickly increases to 100% very consistently across all 10 runs, while test accuracy is also the highest at 78.3%. Other settings that worked well included Comp, DO+TS with a winning test accuracy of 75.8%, and Comp, RS+DO+TS with a winning test accuracy of 77.2%.

(3) **Increasing competitive pressure increases the gap between winning and losing teams:** Another finding is that as more competitive pressure is introduced, the gap between winning and losing teams increases. The winning team increasingly performs better, especially in train accuracy, and the losing team increasingly performs worse than the single team cooperative baseline (i.e. lower than $\sim 50\%$). This confirms our hypothesis that the losing team acts as an *auxiliary* team that boosts the performance of the winning team at its own expense. Furthermore, winning teams that survive through multiple sources of competition learn better communication protocols that allow them to generalize better to new test environments.

(4) **Reward shaping does not improve performance:** In both cooperative and competitive settings, one cannot rely solely on reward shaping to improve generalization. The test accuracies are largely similar across various reward settings.

(5) **Improvement in performance is not due to other factors:** To ensure that the empirical results we observe are not due to confounding factors, we compare the performance in teams of competitive agents with purely cooperative baselines with reward shaping, doubling the number of parameters, and doubling the number of teams. None of these baselines generalize well which shows that the improvement in performance is not due to better rewards, more parameters, or luckier training.

## 5.2 Rates of Convergence

We also compare the convergence in test accuracies of teams trained in both fully cooperative and competitive settings. For test accuracies that ended due to early stopping when train accuracy reached 100%, we propagate the test accuracy corresponding to the best train accuracy over the remaining epochs. This ensures that the test accuracies over multiple runs are averaged accurately across epochs. We outline these results in Figure 3. We find that in all settings, training teams of competitive agents leads to faster convergence and improved performance as compared to the cooperative baselines. Furthermore, composing multiple sources of competition during training steadily improves performance. Figure 3(b) shows a clear trend that: Comp, DO+TS > Comp, DO ≈ Comp, TS > Coop, base. By further adding reward sharing, Figure 3(c) shows that Comp, RS+DO+TS ≈ Comp, RS+DO > Comp, RS+TS > Coop, base.

We find that the fastest convergence in training happens in *earlier stages* of training. Winning teams tend to quickly pull ahead of their losing counterparts during earlier stages and losing teams are unable to recover in later stages of training. This observation is also shared in Table 1 where the gap between winning and losing teams grows as sources of competition are composed. Interestingly, these observations are also mirrored in studies in psychology which argue that competition during childhood is beneficial for early cognitive and social development [7, 37].

Finally, another interesting observation is that high-performing winning teams trained in competitive settings tend to display *lower variance* in test time evaluation as compared to their cooperative counterparts, thereby learning to more stable training.

## 5.3 Emergence of Language

In addition to generalization performance, we also compare the quality of the communication protocols that emerge between the teams of trained agents. Specifically, we measure the signaling that occurs between agents using the Instantaneous Coordination
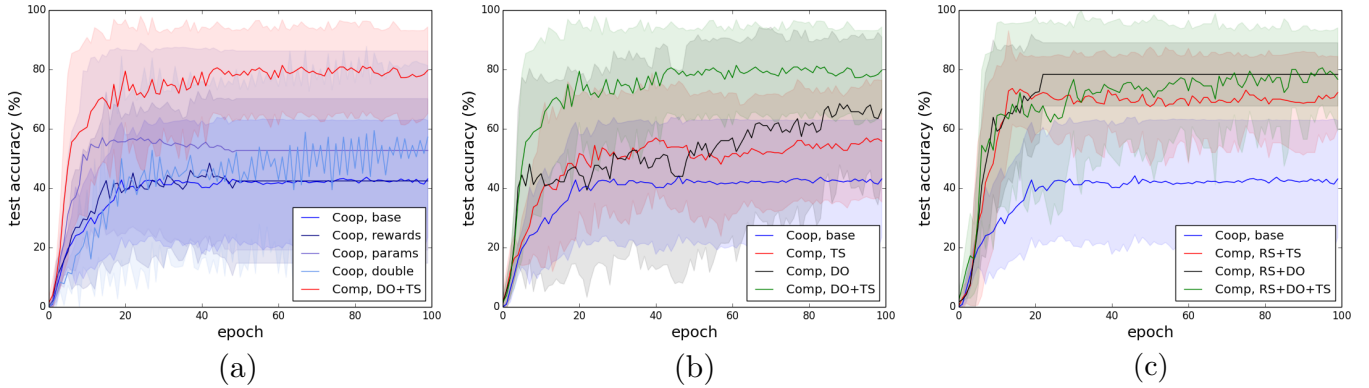
**Figure 3: We plot the convergence of test accuracy across training epochs for the winning team trained in various cooperative and competitive settings. All cooperative baselines are in shades of blue and competitive teams are in red, black, and green. Lines represent the mean across 10 runs and shaded boundaries represent the standard deviations. In (a), we compare the cooperative baselines Coop, base, Coop, rewards, Coop, params, and Coop, double with a well-performing competitive method, Comp, DO+TS. In (b), we compare Coop, base with competitive teams involving Dialog Overhearing (DO) and Task Sharing (TS) (i.e. Comp, TS, Comp, DO, Comp, DO+TS). In (c), we compare Coop, base with competitive teams that additionally incorporate Reward Sharing (RS) (i.e. Comp, RS+TS, Comp, RS+DO, Comp, RS+DO+TS). Winning teams in competitive settings display faster convergence and improved generalization performance in test environments.**

| Type | Method | IC | |
|---|---|---|---|
| | | Winning Team | Losing Team |
| Cooperative baselines | Coop, base [25] | $0.675 \pm 0.099$ | - |
| | Coop, rewards [18, 19] | $0.646 \pm 0.050$ | - |
| | Coop, params | $0.689 \pm 0.101$ | - |
| | Coop, double | $0.719 \pm 0.145$ | $0.691 \pm 0.153$ |
| Competitive methods | Comp, TS | $0.650 \pm 0.139$ | $0.592 \pm 0.128$ |
| | Comp, DO | $0.778 \pm 0.161$ | $0.757 \pm 0.179$ |
| | Comp, DO+TS | $\mathbf{0.806 \pm 0.202}$ | $0.800 \pm 0.204$ |
| | Comp, RS | $0.793 \pm 0.165$ | $0.776 \pm 0.161$ |
| | Comp, RS+TS | $0.726 \pm 0.207$ | $0.718 \pm 0.118$ |
| | Comp, RS+DO | $\mathbf{0.814 \pm 0.154}$ | $0.743 \pm 0.116$ |
| | Comp, RS+DO+TS | $\mathbf{0.834 \pm 0.203}$ | $0.740 \pm 0.142$ |

**Table 2: The Instantaneous Coordination (IC) metric measured across teams trained in various cooperative and competitive settings. All cooperative baselines are in shades of blue and competitive teams are in red. RS: reward sharing, DO: dialog overhearing, TS: task sharing. IC scores are reported separately for winning and losing teams with best accuracies for winning teams in bold. Winning teams in competitive settings perform more informative communication as measured by a higher IC score.**

| Type | Method | $|V_Q|$ | $|V_A|$ | Winning Team |
|---|---|---|---|---|
| Cooperative baselines | Coop, base [25] | 3 | 4 | $45.6 \pm 18.9$ |
| | | 16 | 16 | $26.4 \pm 5.1$ |
| | | 64 | 64 | $22.6 \pm 4.6$ |
| Competitive methods | Comp, RS+DO+TS | 3 | 4 | $\mathbf{77.2 \pm 16.5}$ |
| | | 16 | 16 | $50.8 \pm 26.1$ |
| | | 64 | 64 | $47.5 \pm 25.2$ |

**Table 3: Effect of vocabulary size on both cooperative and competitive training. Similar to Kottur et al. [25], we found that test performance is hurt at large vocab sizes, even under competitive training. For the same fixed vocabulary size, we also see consistent improvements using competitive training as compared to the cooperative baselines, suggesting the utility of our approach across different hyperparameter settings s such as vocabulary sizes.**

(IC) metric [22]. We report these results in Table 2 and focus on the IC between Q-bot and A-bot from the winning team. We observe that IC is highest for the fully competitive setting Comp, RS+DO+TS. Furthermore, by comparing Table 1 with Table 2, we observe a strong correlation between winning teams that signal clearly with high IC scores and winning teams that perform best on test environments. Comp, DO+TS, Comp, RS+DO, and Comp, RS+DO+TS are the training settings that lead to such winning teams. These observations supports our hypothesis that having external

pressure from similar agents encourages the team's Q-bot and A-bot to coordinate better through emergent language, thereby leading to superior task performance which is another benefit of our proposed competitive training method.

Finally, we find that the learned communication protocol is compositional in the same measure as Kottur et al. [25]. For example, Q-bot assigns $Y$ to represent tasks (shape, style), (style, shape), and $X$ for (style, color). The small vocabulary size and memoryless A-bot means that the messages must compose across entities to generalize at test time to unseen instances. We further note that from the convergence graphs as shown in Figure 3, compositionality in emergent language is achieved faster in competitive settings as compared to the fully cooperative counterparts.

We also experimented with large vocab sizes of $|V_Q| = |V_A| = 16$ and 64. We reported these results in Table 3. Similar to Kottur et al. [25], we found that test performance is hurt at large vocab sizes,

| Type | Method | SC (↑) | | H (↓) | |
|---|---|---|---|---|---|
| | | Winning Team | Losing Team | Winning Team | Losing Team |
| Cooperative baselines | Coop, base [25] | $0.631 \pm 0.114$ | - | $1.186 \pm 0.124$ | - |
| | Coop, rewards [18, 19] | $0.640 \pm 0.117$ | - | $1.060 \pm 0.023$ | - |
| | Coop, params | $0.676 \pm 0.132$ | - | $1.138 \pm 0.143$ | - |
| | Coop, double | $0.683 \pm 0.150$ | $0.675 \pm 0.133$ | $1.196 \pm 0.128$ | $1.210 \pm 0.131$ |
| Competitive methods | Comp, TS | $0.610 \pm 0.123$ | $0.618 \pm 0.143$ | $1.089 \pm 0.119$ | $1.107 \pm 0.154$ |
| | Comp, DO | $0.635 \pm 0.149$ | $0.647 \pm 0.185$ | $1.212 \pm 0.123$ | $1.210 \pm 0.135$ |
| | Comp, DO+TS | $0.635 \pm 0.146$ | $0.646 \pm 0.188$ | $1.215 \pm 0.133$ | $1.211 \pm 0.124$ |
| | Comp, RS | $0.677 \pm 0.149$ | $0.658 \pm 0.118$ | $1.205 \pm 0.137$ | $1.207 \pm 0.137$ |
| | Comp, RS+TS | $0.622 \pm 0.157$ | $0.596 \pm 0.143$ | $1.197 \pm 0.143$ | $1.205 \pm 0.130$ |
| | Comp, RS+DO | $\mathbf{0.701 \pm 0.114}$ | $0.670 \pm 0.073$ | $1.174 \pm 0.151$ | $1.217 \pm 0.123$ |
| | Comp, RS+DO+TS | $\mathbf{0.679 \pm 0.219}$ | $0.615 \pm 0.157$ | $1.197 \pm 0.164$ | $1.207 \pm 0.150$ |

Table 4: Speaker Consistency (SC) and Entropy (H) metrics measured across teams in all settings. All cooperative baselines are in shades of blue and competitive teams are in red. RS: reward sharing, DO: dialog overhearing, TS: task sharing. SC and H scores are reported separately for winning and losing teams with best accuracies for winning teams in bold.

even under competitive training. Therefore, we set the vocabulary sizes $|V_Q| = 3$ and $|V_A| = 4$ respectively following Kottur et al. [25]. With these limited vocabulary sizes, we observed good generalization of the language to new object instances. When using large vocabulary sizes, the agents tend to use every vocabulary symbol to memorize pairs of concepts, e.g. symbol $a$ represents a green circle and symbol $b$ represents a green square, etc. instead of representing compositional concepts e.g. symbol $a$ represents the color green and symbol $b$ represents the shape square etc. The compositional vocabulary learned in the latter case is required for generalization to new pairs of concepts at test-time.

From Table 3, it is interesting to note that for the same fixed vocabulary size, we also see consistent improvements using competitive training as compared to the cooperative baselines. This further suggests the utility of our approach across different hyperparameter settings. Moreover, it suggests that competitive training approaches are more robust to different hyperparameter settings such as vocabulary sizes.

### 5.4 Speaker Consistency and Entropy

Here we report the results on two more metrics proposed to measure how informative a language is with respect to the agent's actions [31]: **Speaker Consistency (SC)** measures the mutual information between an agent's message and its future action: $\text{MI}(q_{t(j)}, \hat{w}^G_{i(j)})$ and **Entropy (H)** which measures the entropy of an agent's sequence of outgoing messages. We show these experimental results in Table 4.

In general, competitive teams display a higher speaker consistency score, again showing strong correlation with the best performing teams COMP, RS+DO and COMP, RS+DO+TS. This again implies that the better performing teams trained via competition demonstrate more signaling using their vocabulary. As for entropy, it is hard to interpret this metric [31]. It is traditionally thought that lower entropy in languages represents more compositionality and efficiency in the way meaning is encoded in language. On one hand, it is also possible for an agent to always to send the same symbol which implies the lowest possible entropy, but these messages are unlikely to be informative. The results show that the

entropies across all settings are roughly similar, which we believe imply that the agents are learning communication protocols that are equally complex and rich in nature. However, the improved speaker consistency and instantaneous coordination scores imply that the communication protocols learnt via competition are more informational to the other agents.

## 6 CONCLUSION

In this paper, we revisited emergent language in multi-agent teams from the lens of *competition for performance*: scenarios where competition acts as an additional external pressure for improvement. We start from *Task & Talk*, a previously proposed referential game between two cooperative agents as our testbed and extend it into *Task, Talk & Compete*, a game involving two competitive teams each consisting of cooperative agents. Using our newly proposed *Task, Talk & Compete* benchmark, we showed that competition from an external team acts as social influence that encourages multi-agent populations to develop more informative communication protocols for improved generalization and faster convergence. Our controlled experiments also show that these results are not due to confounding factors such as more parameters, more agents, and reward shaping. This line of work constitutes a step towards studying the emergence of language from agents that are both cooperative and competitive at different levels. Future work can explore the effect of competitive multi-agent training in various real-world settings as well as the emergence of natural language and multimodal dialog.

## 7 ACKNOWLEDGEMENTS

# REFERENCES

[1] Ben Alderson-Day and Charles Fernyhough. 2015. Inner Speech: Development, Cognitive Functions, Phenomenology, and Neurobiology. In *Psychological Bulletin, American Psychological Association.*

[2] Jacob Andreas. 2019. Measuring Compositionality in Representation Learning. In *International Conference on Learning Representations.* https://openreview.net/forum?id=HJz05o0qK7

[3] Bernard J. Baars. 2017. *The Global Workspace Theory of Consciousness.* John Wiley & Sons, Ltd, Chapter 16, 227–242. https://doi.org/10.1002/9781119132363.ch16 arXiv:https://onlinelibrary.wiley.com/doi/pdf/10.1002/9781119132363.ch16

[4] Michael Bacharach. 1991. *Zero-sum Games.* Palgrave Macmillan UK, London, 727–731. https://doi.org/10.1007/978-1-349-21315-3_100

[5] Alan D. Baddeley and Graham Hitch. 1974. Working Memory. Psychology of Learning and Motivation, Vol. 8. Academic Press, 47 – 89. https://doi.org/10.1016/S0079-7421(08)60452-1

[6] Helmy H. Baligh and Leon E. Richartz. 1967. Variable-Sum Game Models of Marketing Problems. *Journal of Marketing Research* 4, 2 (1967), 173–183. https://doi.org/10.1177/002224376700400209 arXiv:https://doi.org/10.1177/002224376700400209

[7] Daphna Bassok. 2012. Competition or Collaboration?: Head Start Enrollment During the Rapid Expansion of State Pre-kindergarten. *Educational Policy* 26, 1 (2012), 96–116. https://doi.org/10.1177/0895904811428973 arXiv:https://doi.org/10.1177/0895904811428973

[8] Ben Bogin, Mor Geva, and Jonathan Berant. 2018. Emergence of Communication in an Interactive World with Consistent Speakers. *CoRR* abs/1809.00549 (2018). arXiv:1809.00549 http://arxiv.org/abs/1809.00549

[9] Antoine Bordes and Jason Weston. 2016. Learning End-to-End Goal-Oriented Dialog. *CoRR* abs/1605.07683 (2016). arXiv:1605.07683 http://arxiv.org/abs/1605.07683

[10] Diane Bouchacourt and Marco Baroni. 2018. How agents see things: On visual representations in an emergent language game. *CoRR* abs/1808.10696 (2018). arXiv:1808.10696 http://arxiv.org/abs/1808.10696

[11] Noam Chomsky. 1957. *Syntactic Structures.* Mouton and Co., The Hague.

[12] F. B. Christiansen and V. Loeschcke. 1990. *Evolution and Competition.* Springer Berlin Heidelberg, Berlin, Heidelberg, 367–394. https://doi.org/10.1007/978-3-642-74474-7_13

[13] Michael Cogswell, Jiasen Lu, Stefan Lee, Devi Parikh, and Dhruv Batra. 2019. Emergence of Compositional Language with Deep Generational Transmission. *CoRR* abs/1904.09067 (2019). arXiv:1904.09067 http://arxiv.org/abs/1904.09067

[14] Abhishek Das, Samyak Datta, Georgia Gkioxari, Stefan Lee, Devi Parikh, and Dhruv Batra. 2018. Embodied Question Answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR).*

[15] Brynne C. DiMenichi and Elizabeth Tricomi. 2015. The power of competition: Effects of social motivation on attention, sustained physical effort, and learning. *Frontiers in Psychology* (2015).

[16] Xavier Glorot, Antoine Bordes, and Yoshua Bengio. 2011. Domain Adaptation for Large-scale Sentiment Classification: A Deep Learning Approach. In *ICML.*

[17] S. Grossberg. 2013. Behavioral economics and neuroeconomics: Cooperation, competition, preference, and decision making. In *The 2013 International Joint Conference on Neural Networks (IJCNN).* 1–5. https://doi.org/10.1109/IJCNN.2013.6706709

[18] Marek Grześ. 2017. Reward Shaping in Episodic Reinforcement Learning. In *AAMAS.*

[19] Marek Grześ and Daniel Kudenko. 2008. Multigrid Reinforcement Learning with Reward Shaping. In *ICANN.*

[20] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation* 9, 8 (1997), 1735–1780.

[21] L. C. Jain and L. R. Medsker. 1999. *Recurrent Neural Networks: Design and Applications* (1st ed.). CRC Press, Inc., Boca Raton, FL, USA.

[22] Natasha Jaques, Angeliki Lazaridou, Edward Hughes, Çaglar Gülçehre, Pedro A. Ortega, DJ Strouse, Joel Z. Leibo, and Nando de Freitas. 2019. Social Influence as Intrinsic Motivation for Multi-Agent Deep Reinforcement Learning. In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA (Proceedings of Machine Learning Research)*, Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.), Vol. 97. PMLR, 3040–3049. http://proceedings.mlr.press/v97/jaques19a.html

[23] Simon Kirby, Hannah Cornish, and Kenny Smith. 2008. Cumulative cultural evolution in the laboratory: An experimental approach to the origins of structure in human language. *Proceedings of the National Academy of Sciences* 105, 31 (2008), 10681–10686. https://doi.org/10.1073/pnas.0707835105 arXiv:https://www.pnas.org/content/105/31/10681.full.pdf

[24] Simon Kirby, Monica Tamariz, Hannah Cornish, and Kenny Smith. 2015. Compression and communication in the cultural evolution of linguistic structure. *Cognition* 141 (2015), 87–102.

[25] Satwik Kottur, José Moura, Stefan Lee, and Dhruv Batra. 2017. Natural Language Does Not Emerge 'Naturally' in Multi-Agent Dialog. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing.* Association for Computational Linguistics, Copenhagen, Denmark, 2962–2967. https://doi.org/10.18653/v1/D17-1321

[26] Angeliki Lazaridou, Karl Moritz Hermann, Karl Tuyls, and Stephen Clark. 2018. Emergence of Linguistic Communication from Referential Games with Symbolic and Pixel Input. In *International Conference on Learning Representations.* https://openreview.net/forum?id=HJGv1Z-AW

[27] Angeliki Lazaridou, Alexander Peysakhovich, and Marco Baroni. 2017. Multi-Agent Cooperation and the Emergence of (Natural) Language. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings.* OpenReview.net. https://openreview.net/forum?id=Hk8N3Sclg

[28] E. G. Leigh Jr. 2010. The evolution of mutualism. *Journal of Evolutionary Biology* 23, 12 (2010), 2507–2528. https://doi.org/10.1111/j.1420-9101.2010.02114.x arXiv:https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.1420-9101.2010.02114.x

[29] Paul Pu Liang, Yao Chong Lim, Yao-Hung Hubert Tsai, Ruslan Salakhutdinov, and Louis-Philippe Morency. 2019. Strong and Simple Baselines for Multimodal Utterance Embeddings. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers).* Association for Computational Linguistics, Minneapolis, Minnesota, 2599–2609. https://doi.org/10.18653/v1/N19-1267

[30] Paul Pu Liang, Ziyin Liu, AmirAli Bagher Zadeh, and Louis-Philippe Morency. 2018. Multimodal Language Analysis with Recurrent Multistage Fusion. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing.* Association for Computational Linguistics, Brussels, Belgium, 150–161. https://doi.org/10.18653/v1/D18-1014

[31] Ryan Lowe, Jakob Foerster, Y-Lan Boureau, Joelle Pineau, and Yann Dauphin. 2019. On the Pitfalls of Measuring Emergent Communication. In *Proceedings of the 18th International Conference on Autonomous Agents and MultiAgent Systems (AAMAS âĂŹ19).* International Foundation for Autonomous Agents and Multiagent Systems, Richland, SC, 693âĂŞ701.

[32] Ryan Lowe, Yi Wu, Aviv Tamar, Jean Harb, Pieter Abbeel, and Igor Mordatch. 2017. Multi-Agent Actor-Critic for Mixed Cooperative-Competitive Environments. In *Proceedings of the 31st International Conference on Neural Information Processing Systems (NIPSâĂŹ17).* Curran Associates Inc., Red Hook, NY, USA, 6382âĂŞ6393.

[33] Dean Mobbs, Demis Hassabis, Rongjun Yu, Carlton Chu, Matthew Rushworth, Erie Boorman, and Tim Dalgleish. 2013. Foraging under Competition: The Neural Basis of Input-Matching in Humans. *Journal of Neuroscience* 33, 23 (2013), 9866–9872. https://doi.org/10.1523/JNEUROSCI.2238-12.2013 arXiv:http://www.jneurosci.org/content/33/23/9866.full.pdf

[34] Martin A. Nowak and David C. Krakauer. 1999. The evolution of language. *Proceedings of the National Academy of Sciences* 96, 14 (1999), 8028–8033. https://doi.org/10.1073/pnas.96.14.8028 arXiv:https://www.pnas.org/content/96/14/8028.full.pdf

[35] Martin A. Nowak, Joshua B. Plotkin, and Vincent A A Jansen. 2000. The evolution of syntactic communication. *Nature* 404 (2000), 495–498.

[36] Minna Pekkonen, Tarmo Ketola, and Jouni T Laakso. 2013. Resource availability and competition shape the evolution of survival and growth ability in a bacterial community. *PLoS One* 8, 9 (2013).

[37] Emmy A. Pepitone. 1985. *Children in Cooperation and Competition.* Springer US, Boston, MA, 17–65. https://doi.org/10.1007/978-1-4899-3650-9_2

[38] Hai Pham, Paul Pu Liang, Thomas Manzini, Louis-Philippe Morency, and Barnabás Póczos. 2019. Found in Translation: Learning Robust Joint Representations by Cyclic Translations between Modalities. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019.* 6892–6899. https://doi.org/10.1609/aaai.v33i01.33016892

[39] Joaquin Quionero-Candela, Masashi Sugiyama, Anton Schwaighofer, and Neil D. Lawrence. 2009. *Dataset Shift in Machine Learning.* The MIT Press.

[40] James H. Read. 2012. Is power zero-sum or variable-sum? Old arguments and new beginnings. *Journal of Political Power* (2012). https://doi.org/10.1080/2158379X.2012.659865 arXiv:https://doi.org/10.1080/2158379X.2012.659865

[41] Kenny Smith, Simon Kirby, and Henry Brighton. 2003. Iterated Learning: A Framework for the Emergence of Language. *Artif. Life* 9, 4 (Sept. 2003), 371–386. https://doi.org/10.1162/106454603322694825

[42] Monica Tamariz and Simon Kirby. 2015. The Cultural Evolution of Language. *Current Opinion in Psychology* 8 (09 2015). https://doi.org/10.1016/j.copsyc.2015.09.003

[43] Paul Vogt. 2005. The emergence of compositional structures in perceptually grounded language games. *Artif. Intell.* 167 (2005), 206–242.

[44] Ronald J. Williams. 1992. Simple Statistical Gradient-Following Algorithms for Connectionist Reinforcement Learning. *Mach. Learn.* 8, 3-4 (May 1992), 229–256. https://doi.org/10.1007/BF00992696

# APPENDIX

## A MODELING THE AGENTS

The *Task, Talk & Compete* game begins with two target instances $I_{(1)}$ and $I_{(2)}$ presented to Q-BOT$_{(1)}$ and Q-BOT$_{(2)}$ respectively, and two tasks $G_{(1)}$ and $G_{(2)}$ presented to A-BOT$_{(1)}$ and A-BOT$_{(2)}$ respectively. Within a team, we largely follow the setting by Kottur et al. [25]. A team consists of agents Q-BOT and A-BOT cooperating in a partially observable world to solve task $G$ given instance $I$. We use lower case characters (e.g. $s_t^Q$) to denote the token symbol and upper case $S_t^Q$ to denote the corresponding representation. We use subscripts to index the rounds and subscripts in parenthesis to index which team the pair of agents belong to (i.e. $s_{t\ (1)}^Q$). We drop the team subscript if it is clear from the context (i.e. same team).

**Base States and Actions:** Each agent observes its specific input (task $G$ for Q-BOT and instance instance $I$ for A-BOT) and the output of the other agent in the same team. At the beginning of round $t$, Q-BOT observes state $s_t^Q = [G, q_1, a_1, \ldots, q_{t-1}, a_{t-1}]$ and utters some token $q_t \in V_Q$. A-BOT observes the history and this new utterance as state $s_t^A = [F, q_1, a_1, \ldots, q_{t-1}, a_{t-1}, q_t]$ and utters $a_t \in V_A$. Note, the two teams share vocabularies (i.e. $V_Q$ is shared by Q-BOT$_{(1)}$ and Q-BOT$_{(2)}$). At the final round, Q-BOT predicts a pair of attribute values $\hat{w}^G = (\hat{w}_1^G, \hat{w}_2^G)$ to solve the task. Q-BOT and A-BOT are modeled as **stochastic policies** $\pi_Q(q_t|s_t^Q; \theta_Q)$ and $\pi_A(a_t|s_t^A; \theta_A)$ implemented as recurrent networks [20, 21]. Q-BOT is modeled with three modules – speaking, listening, and prediction. Given task $G$, Q-BOT stores an initial state $S_{t-1}^Q$ from which it conditionally generates output utterances $q_t \in V_Q$. $S_{t-1}^Q$ is updated using answers $a_t$ from A-BOT and is used to make a prediction $\hat{w}_G$ in the final round. A-BOT is modeled with two modules – speaking and listening. A-BOT encodes instance $I$ into its initial state $S_t^A$ from which it conditionally generates output utterances $a_t \in V_A$. $S_t^A$ is updated using questions $a_t$ from Q-BOT. Q-BOT and A-BOT receive an identical **base reward** of $R$ if Q-BOT's prediction $\hat{w}^G$ matches ground truth $w^G$ and a negative reward of $-10R$ otherwise. $R$ is a hyperparamter which affects the rate of convergence. To train these agents, we update policy parameters $\theta_Q$ and $\theta_A$ using the popular REINFORCE policy gradient algorithm [44].

## B IMPLEMENTATION DETAILS

Table 5 is the list of hyperparameters used. We train all models using the same set of hyperparameters and only modify the rewards and information being shared among agents. All reported results were averaged over 10 runs with randomly initialized random seeds.

In order to report results on the Instantaneous Coordination (IC) metric, we compute the mutual information across the following: within team 1, the mutual information between A-BOT's messages and Q-BOT's first guess, and A-BOT's messages and Q-BOT's second guess. We average these number to obtain an aggregated IC score for team 1 before repeating the procedure for team 2.

| Parameter | Value |
|---|---|
| attribute embeddings size | 20 |
| instance embedding size | 60 |
| $R$ | 100 |
| $|V_Q|$ | 3 |
| $|V_A|$ | 4 |
| $\rho$ | 0.5 |
| batch size | 1000 |
| LSTM dimension | 50 |
| episodes | 1000 |
| max epochs | 50000 |
| learning rate | 0.01 |
| gradient clipping | [-5.0,+5.0] |
| optimizer | Adam |
| num repeats | 10 |

**Table 5: Table of hyperparameters.**