# Exploring an Imagined "We" in Human Collective Hunting: Joint Commitment within Shared Intentionality

**Ning Tang**[1]
ningtangcog@gmail.com

**Siyi Gong**[4]
alicegong@g.ucla.edu

**Minglu Zhao**[3]
minglu.zhao@ucla.edu

**Chenya Gu**[1]
chenyagu@zju.edu.cn

**Jifan Zhou**[1]
Jifanzhou@zju.edu.cn

**Mowei Shen**[1]
mwshen@zju.edu.cn

**Tao Gao**[2,3]
tao.gao@stat.ucla.edu

[1] Department of Psychology, Zhejiang University    [2] Department of Communication, UCLA
[3] Department of Statistics, UCLA    [4] Department of Psychology, UCLA

## Abstract

Human collaboration often involves a decision to pursue one out of multiple comparable goals, in which case it is challenging to remain committed to the same goal collectively. Philosophical theories as well as empirical evidence from developmental psychology suggest that humans, having shared intentionality as an underlying cognitive structure, may be able to form joint commitment in pursuing a collective goal without communication. By conducting experiments in a real-time cooperative hunting game that heavily relies on visual perception, we demonstrated that humans established and maintained robust cooperation with high-quality hunting, even with a large number of potential targets. Additionally, we showed that a Bayesian imagined "We" (IW) model within a joint commitment framework, could capture humans' robustness in resisting alternative targets with relatively high quality of hunting. This poses a contrast with a Reward Sharing (RS) model that, despite performing proficiently in pursuing a single goal, mostly exhibited low-quality hunting and whose teaming fell apart as available targets increased. In a hybrid team simulation experiment, the IW model could better mimic the intentions of human hunters compared to the RS model. Together, the success of the persevered group commitment in humans suggests that shared intentionality is a pivotal element in human cooperation. Moreover, the similarity between the performance of humans and the IW model sheds light on the computational formulation of shared intentionality and further advances our understanding of the nature of cooperation.

**Keywords:** Cooperative hunting; Joint commitment; Shared intentionality; Imagined "We"; Bayesian model

## Introduction

Collective hunting is a complex group activity commonly seen in ecology, in which hunters pursue prey in a group effort. This behavior is evolutionarily significant in the animal kingdom as it extends cooperation beyond kinship to genetically unrelated group mates, including friends, nonfriends (e.g., Seyfarth & Cheney, 2012), and even heterospecifics. For example, Taï Chimpanzees regularly hunt for red colobus monkeys in small groups (e.g., Boesch & Boesch, 1989), and coyotes and badgers were observed to group-hunt ground squirrels (Minta et al., 1992). Moreover, collective hunting and foraging has been viewed as a breakthrough in hominid-evolution and provides a basis for humans' unique large-group cooperation later in phylogenesis (Tomasello, 2014, 2016).

Many differences can be found between human collective hunting and animal collective hunting at the phenomena level, including practices involving sharing spoils, reversing roles, and excluding free riders (Tomasello, 2014). These distinctions are especially salient at the level of commitment. When participating in group activities, it is observed that chimpanzee behavior is primarily motivated by individual desire. An important indicator of this is the absence of obligatory effort to ensure the commitment of other participants in experimental settings (Melis et al., 2006). For example, during collaborative tasks, human-reared chimpanzees did not attempt to re-engage the halfway-quitted experimenters, but instead continued solving the task on their own (Warneken et al., 2006); chimpanzees also stopped acting helpfully once they had received their own shares of reward, in spite of the need for collaboration from their partners (Greenberg et al., 2010). These results are vastly different from human collaborative activities where members are committed to achieving a goal together. Studies show that toddlers as young as 18 months actively attempted to re-engage experimenters who refrained from cooperation, suggesting a desire to regulate group members' commitment (Warneken et al., 2006; Gräfenhain et al., 2009); 3-year-old children continued to collaborate until both partners had received their reward even when the children had received their shares prematurely (Hamann et al., 2012); 3-year-old children also expressed guilt when they accidentally broke a promise and attempted to mend the damage they caused with prosocial, reparative behavior (Vaish et al., 2016).

As shown above, collective hunting and foraging behaviors in humans are more cooperatively structured than those of any other species in the animal kingdom. This behavioral difference implicates deep cognitive roots underlying humans' unique cooperation: They are collaborative interactions in which participants "share" psychological states with one another (Tomasello, 2016). To date, both philosophical theories (e.g., Bratman, 2014; Tuomela, 2007; Searle, 1990; Gilbert, 2013) and empirical research in psychology (e.g., Tomasello, 2014) have pointed out that the species-specific structure at the core of this cognitive representation, is shared intentionality, also known as collective intentionality or joint intentionality.

The notion of collective intentionality—a collective representation of the self and others, or "individuals as a group" (Schweikard & Schmid, 2013)—has been conceptualized throughout history. The idea traces back to early philosophers, such as Aristotle's concept of common striving (*koinonía*) and Rousseau's collective will (*ionvolonté générale*). Sociologists in the last two centuries, such as Durkheim ([1895] 1994) and Weber ([1922] 2009), both have developed their own conceptions of shared intentionality under different labels. More recently, following Collingwood's (1947) definition of "practical social consciousness," Sellars (1974) proposed the concept of "We-intention," emphasizing a non-private attitude involving a shared perspective that can be later used for normative evaluation of contribution. This view has been commonly acknowledged as the predecessor to the modern notion of "collective intentionality" (Schweikard & Schmid, 2013), which was first officially labeled by Searle (1990).

One consensus regarding shared intentionality is that it is irreducible to a sheer summation, aggregation, or distributive pattern of individual intentionality, but rather is a qualitatively different structure of the mind (Schweikard & Schmid, 2013). Among many possible ways to interpret the irreducibility claim, Gilbert focused on what she believed as a definitive

feature of the shared intentionality structure—joint commitment, proposing that a joint intention is only realized when two or more individuals are willing to be "jointly committed to espousing a goal as a body" (Gilbert, 2013). The "goal" here may include a variety of intentions, beliefs, and acceptance (Schweikard & Schmid, 2013), whereas "as a body" indicates an indivisible whole. For example, when individuals A and B jointly believe X as a body, they are committed to forming a single supraindividual agent C that believes X. A and B, in this case of joint commitment, constitute a plural subject that possesses the shared intentional state and cannot be broken down into two single subjects. This notion of plural subject, as we will articulate in detail later, can be understood as a distinct agent with its own actions and mind, including a full set of belief, desire, and intention. Moreover, once the joint commitment to the plural subject is established, any member cannot rescind this commitment unilaterally, and all members in the commitment are normatively responsible for each other. That is, each of them is obligated to act in accordance with their joint goal and entitled to demand others' continuation of the joint action.

Another well-accepted idea, the individual ownership claim, emphasizes that shared intentionality is fundamentally "had by individuals"—it is one's own shared intentionality (Schweikard & Schmid, 2013). This claim clarifies two things about shared intentionality. First, there is no dictatorship—every participant voluntarily forms a joint commitment. Second, shared intention makes sense at the individual level. Instead of blindly executing a subsection of the shared intention that only manifests at a holistic level, each individual owns and understands their intention. These two points correspond to Gilbert's idea that shared intentionality does not require a "single centre of consciousness" or a "distinctive form of 'subjectivity'." Instead, to establish a joint commitment, each individual needs to express a "readiness" for the joint activity, indicating that they are ready and willing to commit, which appeals to all as common knowledge (Gilbert, 2006).

To date, a rich variety of accounts have been proposed regarding how to interpret the "collectiveness" in shared intentionality, as well as how to realize the two claims. For example, alongside Gilbert's focus on joint commitment and the conception of a "body," Bratman (2014) highlights the coordination of plan states of individual agents as the central element in shared intention. Instead of being an intention of a single collective agent, shared intention, on Bratman's account, is an intricate mesh of individual plan-embedded intentions and their interrelations, aligned with each other and commonly known to all. It takes the form of "I intend that we J," where J stands for a joint activity participated in by all intended agents.

Theoretical discussions remark that humans will engage in "whatever behavior" to demonstrate their "readiness" to jointly commit to a goal, as a body, in a shared activity (Gilbert, 2006). As shown in empirical findings, this behavior can be as simple as eye contact(e.g., Siposova et al., 2018). Moreover, based on a rich body of literature on how humans can perceive intentions from motions alone, (e.g., Gao et al., 2009), it is highly likely that the readiness for initiating, as well as maintaining, joint commitment can be expressed simply from coordinated motions. Collectively, these pieces of evidence suggest that humans may be able to establish a sustained joint commitment in a real-time visual-grounded collective hunting task without explicit communication.

So far, current psychophysics works on the perception of intention have been primarily conducted in individual settings and do not concern interactions between multiple agents. For example, a line of studies shows that humans are able to identify the intentions of prey and predators in an online, real-time chasing paradigm (e.g., Gao et al., 2009, 2012; Meyerhoff et

al., 2013). These works heavily revolve around the tension between predators and prey, but rarely involve cooperation between predators. Additionally, participants mostly took the role of observers but not actual players, with a few exceptions where they actually controlled the prey (Gao et al., 2009). Of the few studies that did use displays of cooperative chasing (Yin et al., 2016; Duan et al., 2018), their focus was on perception alone, instead of generating cooperative actions based on perception. Moreover, in these cases, the goal of chasing was fixed to a single target and did not involve the challenge of maintaining joint commitment among many possible goals. These studies generated fruitful results that provide invaluable evidence for humans' ability to infer others' intentions in hunting tasks. However, they cannot be easily generalized to cooperative hunting scenarios in which individuals are not only observers but also participants that generate cooperative behaviors. In such cases, aside from inferring others' intentions and generating action plans accordingly, it is equally important to constantly align one's own intention with others' to converge on a collective goal "as a body" in an enduring fashion. It is thus worth exploring whether humans as engaging players can achieve good cooperation in a real-time hunting task by overcoming such challenges.

Building on currently available studies, we first aim to examine whether a group of three humans can exhibit robust joint commitment while playing a virtual collective hunting game. Following this, we built a computational model of shared intention, named imagined "We," directly inspired by Gilbert's theory of commitment. Our goal is to show that this model can indeed capture important aspects of human cooperative hunting. As a baseline, we also employ a Reward Sharing model without any representation of shared intention. By revealing human performance and comparing it with model performance, we aim to better understand whether joint commitment plays an important role in human cooperation.

## Collective Hunting Experiment

To test the joint commitment in cooperative hunting, we developed a real-time game in a 2-D environment (Fig. 1) with 3 hunters played by humans and 1, 2, or 4 stags as targets which are played by a machine. We want to explore whether human participants consistently converge on the same goal even without communication during the hunting process. Furthermore, we aim to test whether their cooperation can be resistant to an increase in the number of available targets, following the logic that once the joint commitment is achieved, participants should at least be able to secure one target, if not more. Demos of the cooperative hunting process can be found at https://www.youtube.com/playlist?list=PLe7BbCETnjAnCsTkcBEk4k_zuMlpJwZ9i

### Cooperative Hunting Task

Hunters aim to successfully catch the stags, while the stags aim to avoid the hunters. The stags move faster than the hunters, thus requiring hunters to collaborate by persistently chasing a single stag. In multiple-target scenarios, hunters have no predetermined target. Nevertheless, simply for the purpose of improving performance and maximizing accumulated rewards, it is best for them as a group to go after one target at a time. Agents in the environment can take actions within a range of magnitude of force from any direction at a given time step in order to achieve their respective goals. Task performance is evaluated through achieved rewards in a fixed period. The hunters receive a joint reward (+1) upon any successful touch of a stag. The accumulated reward at the end of each trial is used as a dependent measure of the model's performance.
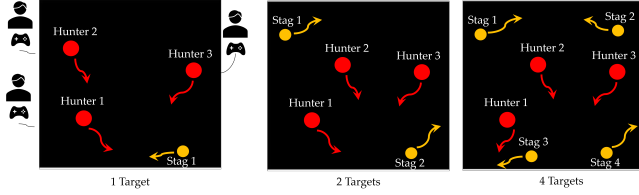
Fig. 1. Cooperative Hunting Task.

## Computational Model

To model human performance, we built an imagined "We" model with a shared intentionality framework to test whether joint commitment is indeed the pivotal mechanism underlying human cooperation.

The previous modeling work on shared intention traces back to early logical models in artificial intelligence, including Grosz and Kraus's "SharedPlans" model (1996) and Levesque et al's definition of a "joint persistent goal" (1990). More recently, some promising modeling works on social agency combined Bayesian inference and Theory of Mind (ToM) approaches with a focus on how agents coordinate their action plans to settle a strategic decision of whether to cooperate or compete in the social world (Kleiman-Weiner et al., 2016; Shum et al., 2019). Under this framework, their formation of cooperation was an abstract planning procedure that showed consistency with Bratman's "meshing of plans" account of shared intentionality. Specifically, the shared intention was implemented as a joint policy that optimizes the team's joint reward, which was then turned into an individual policy by marginalizing the actions of others.

More recently, models of human cooperation integrate the idea of shared intention with normative power. One recent study models multi-agent collaboration in a cooking game using Bayesian Delegation (Wu et al., 2021), in which agents infer what sub-task of a cooking task other agents are working on, and plan accordingly whether they should help with the sub-task or not. In this model, each agent samples a fictitious centralized planner that controls the actions taken by all agents working on the same sub-task. The idea of inferring the states of a fictitious centralized planner resonates with the concept of a shared intention, though the focus here is not on the joint commitment.

**Imagined "We" Model**  Our model builds upon the current progress in Bayesian modeling of shared intention. In addition, our model draws inspiration from Gilbert's plurality subject theory. In our case, human cooperation is assumed, and the focus is on how to model cooperation with a stronger constraint to cohere the team. This is consistent with the perspective that cooperation is qualitatively different from competition (Tomasello, 2009). Preliminary modeling results showing the feasibility of this model with only two collaborators were reported in (Tang et al., 2020) without comparisons to human performance and another baseline model.

Here, we especially focus on the tension between the two well-accepted claims about shared intentionality, being that the collective attitude beyond an individual is, at least to some degree, incompatible with the idea that the intention an individual has cannot escape their own mind. Here, we aim to reconcile this discrepancy by using the imaginative capacity of a causal model (Pearl & Mackenzie, 2018) implemented as the Bayesian Theory of Mind (Baker et al., 2009; Jara-Ettinger et al., 2016). While "We" as a supraindividual agent is not real, each agent can nevertheless imagine the mind of "We" from a collective, "bird's-eye" perspective in their own individual mind (Tomasello, 2009). Specifically, "We" reflects the collective wills of all individuals, but is also a single autonomous agent with its own mind and action just like any
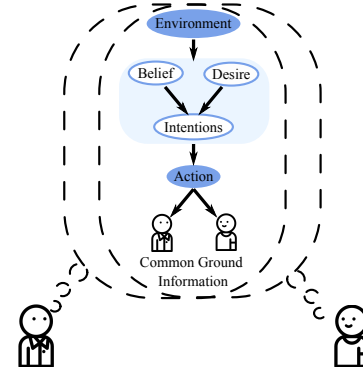


Fig. 2. Imagined "We" Representation. The graphical model in each of the two dashed boxes represents a supraindividual agent "We", which has its own mind containing belief, desire, and intention. Using those mental states, it can rationally control the joint action of the individual agents constituting "We". Each dashed box represents a unique version of this unreal, imagined supraindividual agent "We" inferred by each of the two collaborating individuals here.

ordinary agent, as suggested by Gilbert's theory (2006). Its mental states can be further parsed into belief, desire, and intention which together rationally control this agent's actions. Thus, we can infer the mental states of "We" from its action using ToM, where its state space and action space are simply a concatenation of the state spaces and action spaces of individual agents.

Crucially, this supraindividual agent does not exist in reality—it is ultimately realized by an individual's own imagination about "We" through reasoning counterfactually about "how can an agent explain its own and others' actions if such actions have indeed been rationally controlled by a supraindividual agent 'We'?" For this reason, we call our model the imagined "We" (IW) model (Fig. 2), in part following the classic term of "imagined community," that suggests many communities are first constructed by the imagination (Anderson, 1983). Here we shift the focus from language to perception while perpetuating the same idea: Groups are first imagined before they are formed through practice.

Upon the readiness for joint commitment from all collaborating individuals, each of them, without communication, infers their own version of the imagined "We" by observing the joint action of themselves and their partners in the shared environment. Each agent conceptualizes their own version of "We" and acts by asking "what does 'We' expect me and others to do?" Aside from taking its own action following the intention of "We", an individual agent also expects others to take the actions demanded by "We" (Eq. (1)). Newly generated actions from all agents can be observed and used for each individual agent to update their own inference about "We" for the next time step. An agent's inference is conditioned on the environment in order to capture the intuition that the mind is influenced by the surrounding environment (Eq. (2)). Eventually, joint commitment will be achieved when all individual versions of "We's" are aligned or converged.

$$Joint\ action \sim P(Joint\ action|\text{``We'' } mind, Environment) \tag{1}$$

$$P(\text{``We'' } mind|Joint\ action,\ Environment) \propto$$
$$P(Joint\ action|\text{``We'' } mind,\ Environment) \tag{2}$$
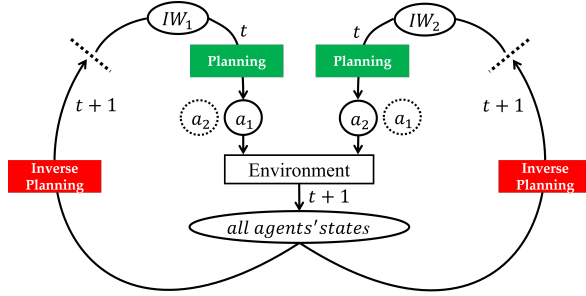$$P(\text{``We'' } mind|Environment)$$

Fig. 3. Bootstrapping Imagined "We". In this case, there are two agents in total, with each loop representing one agent's way to bootstrap an imagined "We". Within an agent's loop, each "IW" node represents a unique distribution of mental states the agent inferred from its imagined "We" agent; the solid "a" node represents the agent's chosen action given its inferred distribution; the dashed "a" node is its expectation of the other agent's action. Actions actually taken by all agents are then observed and used by each agent to update its imagined "We" for the next time step. Here, two agents are used for illustration purposes, while the model can be generalized to multiple agents.

Essentially, this is a process of determining what "We" believes or what "We" wants by observing what "We" has done. Specific to the context of the cooperative hunting task, the environment is fully observable without any uncertainty. The only uncertainty surrounds the intention of "We," concerning "which prey should 'We' pursue persistently?" We model the inference of "We" intention using a bootstrapping method following three steps of computation (Fig. 3).

1. **Goal Sampling**: At each time step ($t$), each agent ($i$) maintains a distribution of intention of "We" as the probability of each target being the joint goal ($GW_{i(t)}$). To decide how to act next, it will draw a sample from this distribution and use it as the estimation for what is the current intention of "We" (Eq. (3)).

$$GW_{i(t)} \sim P(GW_{i(t)}) \qquad (3)$$

2. **Planning**: Given a goal, each agent forms a plan of how "We" should pursue that goal rationally. The output of this planning process is a joint action, including its own action to take, as well as an expectation of the other agents' actions. Thus, each agent is simulating a centralized planner. As an engineering solution, we implemented this rational planning by using a joint policy that was learned through a Multi-Agent Deep Deterministic Policy Gradient (MADDPG, information regarding the MADDPG model will be discussed in detail in the following section) algorithm with only one goal to pursue (Lowe et al., 2017). MADDPG is one of the state-of-the-art implementations within the framework of Multi-Agent Reinforcement Learning (MARL). This joint policy (Eq. (4)) defines the probability of joint actions ($A_{joint(t)} = action_{1(t)}, ..., action_{i(t)}, ..., action_{I(t)}$) conditioning on the current states of "We" ($S_{\text{"We"}(t)} = State_{1(t)}, ..., State_{i(t)}, ..., State_{I(t)}$) and the goal ($S_{GW_{i(t)}}$). Each agent then samples a joint action from the policy distribution and takes its own part of the joint action. Empirically, MADDPG is an algorithm that was found to optimize group reward when only one target was present in a hunting scenario (Zhao et al., 2021). This rational planning phase does not necessarily imply that human cognition employs the exact policy we utilize here. Rather, it is an approxima-

tion of the assumption that humans generally act rationally to optimize their joint utility.

$$P(A_{joint}|GW_{i(t)}) = P(A_{joint}|S_{\text{"we"}, GW_{it}}) \qquad (4)$$

3. **Inference**: After taking one's own action based on the policy determined in the planning phase, each agent observes the actions actually taken by all agents. This enables a Bayesian ToM inference process (Eq. (5)): Conditioning on the observed actions, each hunter computes the posterior probability of a given target being their joint goal. After updating the posterior of the Imagined "We" mind, each agent goes back to step 1, sampling a new goal and repeating the process.

$$P(GW_{i(t+1)}|A_{joint(t)}) \propto P(A_{joint(t)}|GW_{i(t)})P(GW_{i(t)}) \qquad (5)$$

**Baseline Model**  To explore the necessity of modeling a shared intention framework, we additionally examine cooperation in a Reward Sharing (RS) model as a baseline. Here we use the MADDPG, an algorithm within the MARL framework. MARL adopts the perspective that social skills are learned through trial-and-error (Hayek, 2011). It has been successfully applied to complex multi-agent coordination tasks (Berner et al., 2019; Vinyals et al., 2019), in which it splits group rewards evenly among all agents by assigning them the same reward function. Since we used the MADDPG algorithm to train the joint policy in the one-target scenario, the RS model and the IW model are identical in the case of a single target. For the multiple target conditions, we further trained a separate RS model for each set size of targets. Note that this separate training is not required in the IW model as the inference of a goal is achieved by Bayesian inference and the same one-target policy is applied to all conditions for joint goal inference and planning.

Despite being a reinforcement model, the RS model displays several interesting components that can be considered precursors of ToM. For example, it acts based on predicting what other agents will do given their current states. Then it evaluates the utility of its own action given the current joint state of all agents plus its prediction of other agents' actions. As a type of reinforcement learning model, it has a generic framework that can be universally applied to any multi-agent scenario, including both cooperation and competition, which only differ by whether agents' rewards are aligned or opposed. At its core, it is the opposite of the IW model that assumes cooperation is qualitatively different from competition and thus requires a brand new cognitive scaffold. In short, collaborations in the RS model are encouraged by sharing rewards without any reference to commitment, whereas teaming behaviors in the IW model are enforced by shared intention.

**Model Task & Prediction**  The same task completed by human participants was used to test the IW model and the RS model and compare their performance to that of humans. We aim to explore whether human performance in cooperative hunting can be better captured by the stronger-constrained IW model or the weaker-constrained MARL model.

## Human Experiment

Thirty-three (3 participants in 1 group, 11 groups total) students (14 females, 19 males) participated in this experiment. All were between the ages of 21 and 28 ($M_{age} = 23$, $SD = 2.0$) with normal or corrected-to-normal visual acuity. All participants signed the informed consent form and received experimental rewards related to performance after the experiment
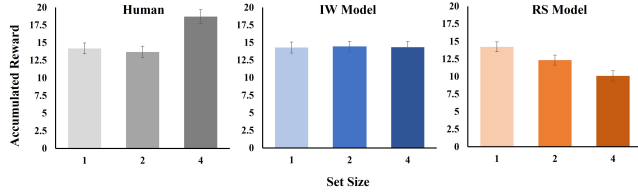
Fig. 4. Results of the accumulated rewards.



Fig. 5. Results of the percentages of different quality rewards.



Fig. 6. Results of the entropy of touched target distribution.

ended.

**Material & Procedure** 20 trials were set up for each condition with different set sizes of hunting. The hunting task was presented on a $38.0° \times 38.0°$ window displayed simultaneously on three monitors, one for each participant. Each participant was instructed to use an Xbox controller to control the simulated physical forces they could apply to drive the hunters on the screen. Three hunters were represented by circular shapes with a diameter of $1.9°$ and colors red (205, 0, 0), green (0, 139, 0), and blue (0, 0, 205). The stags are circular shapes with a diameter of $1.3°$ and could be distinguished by their colors (255, 165, 0), (255, 127, 36), (255, 165, 79), and (255, 193, 37).

**Results  Overall Performance.** We first analyzed the accumulated reward of the hunting game, which reflected the overall performance of humans and models (Fig. 4). One-way ANOVA revealed a significant main effect of set size ($F(2, 30) = 6.58$, $p < .01$, $\eta_p^2 = 0.31$). The post-hoc comparisons showed that when the number of targets increased to 4, the performance was even higher than the 1 and 2 target conditions ($ps < .01$). Overall, when the number of targets increased, the accumulated reward of human hunters did not decrease but increased instead. For the IW model, the main effect of set size was not significant ($F(2, 30) = 1.15$, $p = .331$, $\eta_p^2 = .07$). The results revealed that the performance of the IW model did not decrease as the number of targets increased. For the RS model, a significant main effect of set size was revealed ($F(2, 30) = 87.79$, $p < .01$, $\eta_p^2 = .86$). The post-hoc comparisons showed that when the number of targets increased, the performance gradually decreased ($ps < .05$ for both set size 1-2 and set size 2-4 comparisons).

**Quality of Hunting.** Besides the above quantitative analysis of the overall performance, we further explore the quality of hunting in humans and models (Fig. 5). Here we indicate the quality of hunting by measuring the "duration of touch," defined by the number of consecutive time steps in which at least one hunter touches the stag. In real life, a short touch duration suggests a hit to the target, but not necessarily a catch, whereas a long touch duration indicates a greater likelihood for a real catch or kill, as hunter(s) may have cornered the stag. Thus, the quality of raw rewards was categorized into 3 classes in terms of touch duration: low (1 time step), median (2 time steps), or high (3 or more time steps). The percentages of different qualities of rewards in the total rewards were then measured. Note that this is a post-hoc analysis of hunting quality based on previous results. The quality of hunting was not part of the instructions for both humans and models, and thus neither were optimizing their performance on this metric. The fact that without any instruction, the quality of human hunting was higher than those of models is further discussed in the Discussion.

For the set size 2 condition, one-way ANOVA revealed a significant main effect of player type in both low- and high-quality conditions ($F(2, 30) = 523.69$, $p < .001$, $\eta_p^2 = 0.97$; $F(2, 30) = 672.49$, $p < .001$, $\eta_p^2 = 0.98$). The post-hoc comparisons showed that humans received low-quality reward
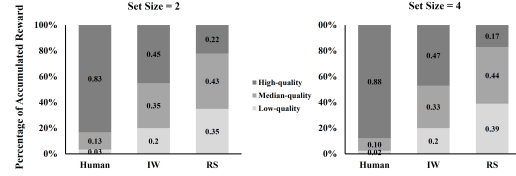
less often than the IW model ($p < .001$), which received it less often than the RS model ($ps < .001$). On the contrary, humans received high-quality reward more often than the IW model ($p < .001$), which received it more often than the RS model ($p < .001$). Similar main effects of player type ($F(2, 30) = 576.61$, $p < .001$, $\eta_p^2 = 0.97$; $F(2, 30) = 803.57$, $p < .001$, $\eta_p^2 = 0.98$) were found in set size 4 condition with similar post-hoc comparisons ($ps < .001$). These results collectively suggest that human hunters achieved the largest proportion of high-quality hunting, followed by the IW hunters, while the RS hunters achieved the smallest proportion. Notably, the IW model had a relatively high quality of hunting, though there was still room for improvement.

**Goal Consistency.** Beyond task performance, we further analyzed the goal consistency among hunters in a team (Fig. 6). Here we measured the entropy of the distribution of the touched target, of which a lower entropy value indicates a higher convergence or concentration on the same goal from all hunters. Both set size 2 and 4 conditions showed a significant main effect of agent type ($F(2, 30) = 110.11$, $p < .01$, $\eta_p^2 = 0.87$; $F(2, 30) = 13.38$, $p < .001$, $\eta_p^2 = 0.44$). For the set size 2 condition, the entropy of the touched target distribution of humans was higher than that of the IW model ($p < .001$), but was lower than that of the RS model ($p < .001$). For the set size 4 condition, the difference between the entropy of humans and that of the IW model was not significant ($p = .24$), but both of them were higher than the entropy of the RS model ($ps < .01$). These results suggested that the way humans pursued their goals could be better captured by the IW model than the RS model.

**Hybrid Team Simulation**

**Overview** Thus far we have only examined the performance of each type of player within their homogeneous group. Here, we take one step further to measure how well they can cooperate with each other to investigate the compatibility between their hunting strategies. We conducted a hybrid team simulation experiment based on the pre-recorded trajectories of all agents in the human experiment. To see how well they could cooperate with each other, we replaced a human hunter with an IW or RS model hunter while leaving the trajectories of all other agents untouched. As the new hunter was "invisible" to the pre-recorded stags, we expected an overall increase in the performance of the hybrid team compared to the original all-human team, but we were more interested to discover whether the models could align their goals with the rest of the human hunters. We examined the matching
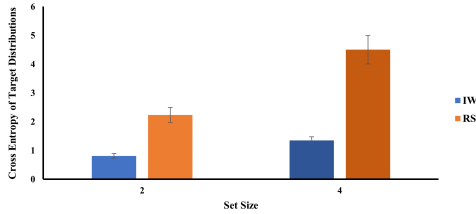
Fig. 7. Results of the cross entropy of the touched target distributions between the hybrid teams and the original all-human team.

of goal consistency by measuring the cross entropy of the touched target distributions between the hybrid teams and the original all-human team. If the models can successfully infer the goal of human hunters and cooperate by coordinating their hunters' behaviors to commit to the same goal, then the touched target distribution of the hybrid team would be similar to that of the all-human team with a low cross entropy.

**Results**  For both set size 2 and 4 conditions (Fig. 7), the cross entropy between the IW-human team and all-human team was significantly lower than that between the RS-human team and all-human team. ($t(20)$ = -17.13, $p < .001$, $d$ = 7.67; $t(20)$ = -37.22, $p < .001$, $d$ = 16.65). This shows that the IW model could better "replicate" the intention of the human player they replaced, than the RS model.

## Discussion

The human behavioral results demonstrate a successful expansion of human social perception from individual chasing tasks to a multi-agent cooperation task involving the integration of perception and planning. We thus show that humans were indeed capable of, and in fact good at, achieving effective collaboration in a hunting task with multiple temptations. Moreover, this high-level performance in humans was achieved without any form of explicit communication. In fact, the only perceptual inputs were the dots moving on the display. This result corroborates the theory that communication may have only emerged in environments where collaboration already existed (Tomasello, 2010). Although we did not directly model communication, our results echo the notion that the more we can achieve without communicating, the more effective communication will be when we do communicate; thus supporting that communication can be achieved with only sparse, highly context-dependent input. This idea is also highlighted by a recent study showing that a "We" built by visual common sense takes on much of the heavy-lifting in communication, thus enabling humans to produce and understand indirect and ambiguous signals in and from few words (Stacy et al., 2020).

From the overall performance and the entropy analysis, we saw that cooperation in humans and the IW model were reasonably robust—their teams were consistent in goal pursuit and their accumulated rewards were well maintained, even when the number of potential targets escalated. This result is non-trivial as evidenced by the RS model whose performance plummeted with increased potential targets, despite the model's specific training for a large set size. The similarity in performance between humans and the IW model was further confirmed in the hybrid team simulation, in which the IW model aligned its goal with human hunters relatively well. This result is especially noteworthy as the IW model was purely trained from the MADDPG model combined with Bayesian inference, but never from human trajectories. The successful cooperation between humans and the IW model indicates that both of them are better at coordinating their actions for achieving a joint goal than the RS model, thus

supporting the claim that shared intentionality is a key mechanism in enabling humans to stay robustly committed in cooperation. Moreover, even though success was defined by the mere touch of a stag by any hunter, the majority of the human hunts were of high quality, a characteristic partially reflected in the IW model hunts. This phenomenon reveals that, to at least some degree, there is a spontaneous, agreed-upon emphasis on the quality of cooperation in both humans and the IW model. Cognitively, this may reflect an important strategic and functional aspect of shared intentionality in the context of coordination—in real life, cornering a stag is a much more effective and sophisticated coordination of action to ensure its capture as compared to a hit. Future behavioral paradigms could emphasize high-quality kills over low-quality hits, for calibrating successful cooperation.

In contrast to the IW model, overall the RS model performed less proficiently as the number of targets increased, showed much lower goal consistency between team members, exhibited a lower level of compatibility when teaming with humans, and made more frequent touches of lower quality. These results should not be taken lightly, as the MADDPG algorithm is one of the best models for coordinating multi-agent chasing when there is only one target—it is why we employed it as the base model for the IW model. Nevertheless, this model failed to handle commitment when facing a large number of targets, which happened to be the most important aspect of cooperation considered by Gilbert (2006). These results suggest that as collaboration tasks in the real world are often tempted by many different desires, sharing rewards only provides a weak constraint on collective behavior. Related to this challenge, a recent study has demonstrated MADDPG's difficulty in handling the free-rider problem—another major challenge in cooperation (Zhao et al., 2021). The difference in how the RS model pursued goals as compared to humans is also consistent with a recent study on human-AI teaming (Siu et al., 2021). The RL model, based on pure learning, received lower subjective scores on human-rated performance, teamwork, interpretability, and trust as compared to other models, though the objective performance of the human-AI team was identical to other models.

It is nevertheless also true that as compared to the IW model, humans accumulated more rewards when the set size increased to 4, achieved high-quality hunts more often, and showed greater flexibility in goal pursuit in the set size 2 condition, revealing that humans exceeded the IW model in cooperative hunting in both quantity and quality. This reflects greater flexibility in human cooperative behaviors and may even reveal other higher-level aspects of shared intentionality in humans, such as the malleable nature of joint commitment. For example, human participants, while committed to the collective superordinate intention, might have also believed that it was not necessary nor efficient to dispatch all three hunters to one target in cases where two hunters already had control of the prey. The advantage of having this flexibility is especially conspicuous when there are larger numbers of prey to hunt, providing more opportunities for humans to efficiently allocate resources to obtain greater rewards, as evident by humans' superior performance in the set size 4 condition over the set size 2 condition. This performance difference between conditions, however, is not seen in the IW model, likely due to its rigidity in team structures—its definition of commitment always assumed the 3 hunters as a group and thus neglected other potential structurings that allowed for more fruitful reward. We believe that addressing the challenge of how to integrate flexible task assignments while maintaining the constraint of shared intentionality, will be the next step in advancing cooperation modeling in the future.

# References

Anderson, B. (1983). *Imagined communities: Reflections on the origin and spread of nationalism*. Verso books.

Baker, C. L., Saxe, R., & Tenenbaum, J. B. (2009). Action understanding as inverse planning. *Cognition*, *113*(3), 329-349.

Berner, C., Brockman, G., Chan, B., Cheung, V., Debiak, P., Dennison, C., ... others (2019). Dota 2 with large scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680*.

Boesch, C., & Boesch, H. (1989). Hunting behavior of wild chimpanzees in the tai national park. *American journal of physical anthropology*, *78*(4), 547–573.

Bratman, M. (2014). Rational and social agency: Reflections and replies. *Rational and Social Agency*, 294–343.

Collingwood, R. G. (1947). *The new leviathan: Or man, society, civilization and barbarism*. Clarendon Press.

Duan, J., Yang, Z., He, X., Shao, M., & Yin, J. (2018). Automatic attribution of social coordination information to chasing scenes: evidence from mu suppression. *Experimental brain research*, *236*(1), 117-127.

Durkheim, E. ([1895] 1994). On social facts. In M. Micheal & M. L. C. (Eds.), *Readings in the philosophy of social science*. MIT Press.

Gao, T., Newman, G. E., & Scholl, B. J. (2009). The psychophysics of chasing: A case study in the perception of animacy. *Cognitive psychology*, *59*(2), 154–179.

Gao, T., Scholl, B. J., & McCarthy, G. (2012). Dissociating the detection of intentionality from animacy in the right posterior superior temporal sulcus. *Journal of Neuroscience*, *32*(41), 14276–14280.

Gilbert, M. (2006). *A theory of political obligation: Membership, commitment, and the bonds of society*. Oxford University Press.

Gilbert, M. (2013). *Joint commitment: How we make the social world*. Oxford University Press.

Gräfenhain, M., Behne, T., Carpenter, M., & Tomasello, M. (2009). Young children's understanding of joint commitments. *Developmental psychology*, *45*(5), 1430–1443.

Greenberg, J. R., Hamann, K., Warneken, F., & Tomasello, M. (2010). Chimpanzee helping in collaborative and noncollaborative contexts. *Animal Behaviour*, *80*(5), 873–880.

Grosz, B. J., & Kraus, S. (1996). Collaborative plans for complex group action. *Artificial Intelligence*, *86*(2), 269–357.

Hamann, K., Warneken, F., & Tomasello, M. (2012). Children's developing commitments to joint goals. *Child development*, *83*(1), 137–145.

Hayek, F. A. (2011). *The constitution of liberty* (5th ed.; R. Hamowy, Ed.). University of Chicago Press.

Jara-Ettinger, J., Gweon, H., Schulz, L. E., & Tenenbaum, J. B. (2016). The naïve utility calculus: Computational principles underlying commonsense psychology. *Trends in cognitive sciences*, *20*(8), 589-604.

Kleiman-Weiner, M., Ho, M. K., Austerweil, J. L., Littman, M. L., & Tenenbaum, J. B. (2016). Coordinate to cooperate or compete: abstract goals and joint intentions in social interaction. In A. Papafragou, D. Grodner, D. Mirman, & J. C. Trueswell (Eds.), *Proceedings of the 38th annual conference of the cognitive science society* (p. 1679-1684). Austin, TX: Cognitive Science Society.

Levesque, H. J., Cohen, P. R., & Nunes, J. H. (1990). On acting together. In *Proceedings of the eighth aaai conference on artificial intelligence*. American Association for Artificial Intelligence.

Lowe, R., Wu, Y., Tamar, A., Harb, J., Abbeel, P., & Mordatch, I. (2017). Multi-agent actor-critic for mixed cooperative competitive environments. *arXiv preprint arXiv:1706.02275*.

Melis, A. P., Hare, B., & Tomasello, M. (2006). Engineering cooperation in chimpanzees: tolerance constraints on cooperation. *Animal Behaviour*, *72*(2), 275-286.

Meyerhoff, H. S., Huff, M., & Schwan, S. (2013). Linking perceptual animacy to attention: Evidence from the chasing detection paradigm. *Journal of Experimental Psychology: Human Perception and Performance*, *39*(4), 1003-1015.

Minta, S. C., Minta, K. A., & Lott, D. F. (1992). Hunting associations between badgers (taxidea taxus) and coyotes (canis latrans). *Journal of Mammalogy*, *73*(4), 814–820.

Pearl, J., & Mackenzie, D. (2018). *The book of why: the new science of cause and effect*. Basic books.

Schweikard, D. P., & Schmid, H. B. (2013). Collective intentionality. In E. N. Zalta (Ed.), *The stanford encyclopedia of philosophy* (Fall 2021 ed.). Metaphysics Research Lab, Stanford University.

Searle, J. R. (1990). Collective intentions and actions. In P. Cohen, J. Morgan, & M. Pollack (Eds.), *Intentions in communication*. MIT Press.

Sellars, W. (1974). *Essays in philosophy and its history*. Springer Science Business Media.

Seyfarth, R. M., & Cheney, D. L. (2012). The evolutionary origins of friendship. *Annual review of psychology*, *63*, 153–177.

Shum, M., Kleiman-Weiner, M., Littman, M. L., & Tenenbaum, J. B. (2019). Theory of minds: Understanding behavior in groups through inverse planning. In *Proceedings of the aaai conference on artificial intelligence* (Vol. 33, p. 6163-6170).

Siposova, B., Tomasello, M., & Carpenter, M. (2018). Communicative eye contact signals a commitment to cooperate for young children. *Cognition*, *179*, 192–201.

Siu, H. C., Peña, J., Chen, E., Zhou, Y., Lopez, V., Palko, K., ... Allen, R. (2021). Evaluation of human-ai teams for learned and rule-based agents in hanabi. *Advances in Neural Information Processing Systems*, *34*.

Stacy, S., Zhao, Q., Zhao, M., Kleiman-Weiner, M., & Gao, T. (2020). Intuitive signaling through an" imagined we'". In B. Armstrong, S. Denison, S. Mack, & Y. Xu (Eds.), *Proceedings of the 42st annual conference of the cognitive science society* (p. 1880-1881). Cognitive Science Society.

Tang, N., Stacy, S., Zhao, M., Marquez, G., & Gao, T. (2020). Bootstrapping an imagined we for cooperation. In B. Armstrong, S. Denison, S. Mack, & Y. Xu (Eds.), *Proceedings of the 42st annual conference of the cognitive science society* (p. 2453-2458). Cognitive Science Society.

Tomasello, M. (2009). *Why we cooperate*. MIT press.

Tomasello, M. (2010). *Origins of human communication*. MIT press.

Tomasello, M. (2014). The ultra-social animal. *European journal of social psychology*, *44*(3), 187–194.

Tomasello, M. (2016). *A natural history of human morality*. Harvard University Press.

Tuomela, R. (2007). *The philosophy of sociality: The shared point of view*. Oxford University Press.

Vaish, A., Carpenter, M., & Tomasello, M. (2016). The early emergence of guilt-motivated prosocial behavior. *Child Development*, *87*(6), 1772–1782.

Vinyals, O., Babuschkin, I., Czarnecki, W. M., Mathieu, M., Dudzik, A., Chung, J., ... others (2019). Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature*, *575*(7782), 350–354.

Warneken, F., Chen, F., & Tomasello, M. (2006). Cooperative activities in young children and chimpanzees. *Child development*, *77*(3), 640–663.

Weber, M. ([1922] 2009). *The theory of social and economic organization*. Simon and Schuster.

Wu, S. A., Wang, R. E., Evans, J. A., Tenenbaum, J. B., Parkes, D. C., & Kleiman-Weiner, M. (2021). Too many cooks: Bayesian inference for coordinating multi-agent collaboration. *Topics in Cognitive Science*, *13*(2), 414–432.

Yin, J., Xu, H., Ding, X., Liang, J., Shui, R., & Shen, M. (2016). Social constraints from an observer's perspective: Coordinated actions make an agent's position more predictable. *Cognition*, *151*, 10-17.

Zhao, M., Tang, N., Dahmani, A. L., Perry, R. R., Zhu, Y., Rossano, F., & Gao, T. (2021). Sharing is not needed: Modeling animal coordinated hunting with reinforcement learning. In T. Fitch, C. Lamm, H. Leder, & K. Teßmar-Raible (Eds.), *Proceedings of the 43st annual conference of the cognitive science society* (p. 1014-1015).