

UC Merced

Proceedings of the Annual Meeting of the Cognitive Science Society

Title

Moment-to-moment decisions of when and how to help another person

Permalink

<https://escholarship.org/uc/item/5nf3233k>

Journal

Proceedings of the Annual Meeting of the Cognitive Science Society, 46(0)

Authors

Osborn Popp, Pamela Joy

Gureckis, Todd M

Publication Date

2024

Peer reviewed

Moment-to-moment decisions of when and how to help another person

Pamela J. Osborn Popp (pamop@nyu.edu)

Center for Neural Science, New York University

Todd M. Gureckis (todd.gureckis@nyu.edu)

Department of Psychology, New York University
New York, NY 10003 USA

Abstract

Helping is a universal human behavior, and is a core aspect of a functioning society. However, the decision to provide help, and what type of help to provide, is a complex cognitive calculation that weights many costs and benefits simultaneously. In this paper, we explore how various costs influence the moment-to-moment decision to help in a simple video game. Participants were paired with another human participant and were asked to make repeated decisions that could benefit either themselves or their partner. Several preregistered manipulations altered the cost each person paid for actions in the environment, the intrinsic resource capacity of individuals to perform the task, the visibility of the other player's score, and the affordances within the environment for helping. The results give novel insight into the cost-benefit analyses that people apply when providing help, and highlight the role of reciprocity in influencing helping decisions.

Keywords: decision making; social cognition; helping; collaboration; altruism; reciprocity

Introduction

Helping others is central to our lives but also mysterious from a computational perspective. A large body of research in economics and psychology has described helping as an altruistic behavior where an agent sacrifices local personal gains for less tangible rewards of self-presentation, future reciprocity, and collective benefit (Marsh, 2016; Andreoni & Miller, 2008; Fehr & Fischbacher, 2003; Tomasello, 2009). However, the local decision to help is often made quickly in the moment and simultaneously balances many complex costs and considerations. Perhaps most importantly the decision of *if* to help depends on an analysis of what help is needed. For example, if a person stops you on the street to ask for directions you might be happy to offer verbal directions, but if they expected you to travel several miles with them, or carry them to their destination, you might be less willing to assist. You might also consider several other factors like your personal suitability to help in comparison to the other person (e.g., we understand children to be more limited than adults and thus help them in cases where it would seem odd to help an adult). In the span of a short interaction, even with a stranger, we seem to effortlessly analyze these and other factors to determine when we should help and how.

The underlying computations supporting how people understand when and how to help another person remain unclear. One of the biggest challenges when analyzing social decision-making is the vast increase in sources of uncertainty

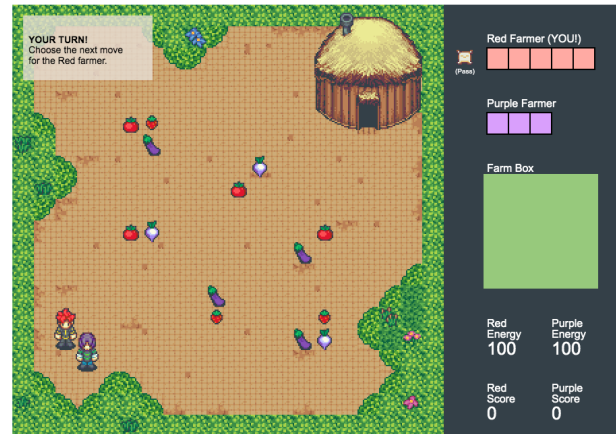


Figure 1: Experiment display from the perspective of the participant controlling the Red farmer character.

compared to individual decision-making tasks (Ho, MacGlashan, Littman, & Cushman, 2017; Kleiman-Weiner, Ho, Austerweil, Littman, & Tenenbaum, 2016). In helping, these include what help another person needs, what their goals are, and the opportunity cost for providing help, among many other factors. Traditional approaches within economics and behavioral decision making attempt to simplify these concerns using simple tasks where altruistic decisions to help are obvious and discrete (see, e.g., Boyd (1988); Trivers (1971)). However, in more realistic environments, individual helping actions are complex, do not involve fixed or known costs and involve sequences of actions that unfold over time. For example, helping someone cross the street means helping them navigate all the potential road hazards and may take a variable amount of time and energy. Unraveling these nuances is critical to building a more comprehensive cognitive account of human helping and collaboration.

Overview of study

In the present study we designed a novel experimental task to examine how people decide if and how to help¹. The task

¹ It is helpful to clarify the distinction between helping and collaboration. One unique feature of helping is that the two agents have individual goals (separate reward functions). When one agent helps another, they devote effort towards the partner agent's goal, potentially sacrificing their own reward or incurring costs to do so. In

was structured such that each player had their own distinct rewards, but also the ability to affect the other player's rewards. More specifically, in the task, participants played as a "red farmer" and a "purple farmer" tasked with picking up vegetables and delivering them to a barn (see Fig. 1). Vegetables were either red or purple, with red vegetables contributing exclusively to the red farmer's reward and purple vegetables contributing exclusively to the purple farmer's reward. Players could pick up either color of vegetable, and were incentivized to perform the harvest efficiently because walking around the farm was costly.

Under this design, participants' actions can be characterized by whether the action mainly benefits them (e.g., picking up their own color vegetable) or their partner (e.g., picking up their partner's color vegetable). At the same time, the decision of *how* to help is quite undefined. While help always involves picking up the other agent's vegetables, which vegetables to pick up matters quite substantially (e.g., picking up an object very close to the other agent might force them to walk even further for their next item). Additionally, different amounts of energy can be expended on actions depending on the walking distance from the current location. Helping might also unfold in real time across many turns, where help is provided or withheld based on an ongoing assessment of the local context. We define different measures of helping that account for the rate at which players choose helpful actions, as well as how costly those actions are, and how they vary based on the features of the environment and partner.

We predicted that participants would help more when task costs were lower, and when they could see their partner's score as well as their own (increasing awareness of the other person's situation). We also predicted that participants would help more when they had greater abilities than their partner (specifically, a larger capacity backpack in which to carry vegetables). Finally, we predicted that the affordances of the environment – how different subtasks were distributed spatially – would influence helping behavior. Under this view, helping is "opportunistic" and occurs at moments where the energy expenditure to help is momentarily lowered. We tested these predictions through pre-registered analyses of our design.

Critically, in none of the conditions were subjects required to help (and, in fact, the reward structure of the task discouraged helping). However, we expected spontaneous and natural examples of helping to emerge based on the hypothesized variables above as well as others including expectations of reciprocity that emerge over repeated interactions (Stephens, 1996; Fehr & Gächter, 2000).

Experimental method – The Farm Task

The experiment was designed to create a small "virtual world" where individual agents pursue their own rewards but have the possibility of helping other agents in the same en-

contrast, collaborative agents generally work to maximize a common goal shared by all players.

vironment. The experiment was conducted in a real time, turn-based video game with a randomly paired human partner (see Fig. 1). The experimental protocol and preliminary analyses were pre-registered², and a full demo of the game, as well as gameplay from every recorded participant session, is available on a project website³.

Participants

Participants were recruited on Prolific (<https://www.prolific.co>; Palan and Schitter (2018)) to take part in a psychology experiment. Subjects were informed that the task would require them to play a video game with another online player set around collecting the harvest on a farm. After passing a comprehension check of the instructions and completing a brief CAPTCHA (Von Ahn, Blum, Hopper, and Langford (2003)) task, participants entered a waiting room to be paired with another player online. The task began immediately after participants were assigned to a pair; if participants did not find a pair within five minutes, they exited the experiment and received partial compensation for their time at a rate of \$15 per hour. The full experiment took about 30 minutes and participants were paid a \$7.50 base rate plus a bonus of up to \$5 depending on their performance (an average pay rate exceeding the local minimum wage where data were collected).

We collected data in March and April of 2023 from 750 participants who were paired and began the game together (375 dyads). Dyads were excluded from analysis and modeling if their data was incomplete, e.g., if one participant left the experiment early. The final dataset includes complete data from $N = 628$ participants (314 dyads).

Experiment Design

In the experiment, two anonymous Prolific users were randomly paired together to play 12 rounds ("games") of the farm task in one sitting ("session"). In each game, participants controlled small avatars on a virtual two-dimensional "farm" (similar to well known video games like *Harvest Moon* or *Stardew Valley*). The farm was composed of an open field with several recognizable objects (vegetables) arranged in various locations on a grid. The objects were either red (strawberries or tomatoes) or purple (eggplant or turnips). Players were randomly assigned at the start of the task to control either the "red" farmer avatar or the "purple" farmer avatar (lower left corner of Fig. 1). The goal of both players was to efficiently pick up and deliver the vegetables to a barn where objects could be stored (upper right corner of the farm in Fig. 1). A game ended once all vegetables were successfully delivered to the barn. Importantly, nowhere in the instructions was it suggested that one player might help the other.

Players began each round with a set amount of "energy," which was depleted proportional to the distance the avatar "walked" to encourage efficient strategies for collecting the

²<https://aspredicted.org/QGG-SKJ>

³<https://exps.gureckislab.org/e/helping-game-viewer/>

harvest. Critically, players only earned bonus points as a function of the quantity of their own color vegetables delivered and their own remaining energy. For instance, the red farmer’s score was computed as the number of red vegetables multiplied by the red farmer’s remaining energy units at the end of the round. However, players could choose to pick up vegetables of either color, choosing to incur energy costs to benefit their partner (i.e., “helping”). Participant might help one another due to altruism, a desire to complete the game faster, or any number of alternative reasons, which we revisit later in the discussion. As part of the design, players had a finite backpack capacity which limited how many vegetables they could carry at one time before they had to drop off their harvest at the barn.

Within each game, one participant was selected randomly to go first. Then, participants alternated turns in which they clicked on a target object to direct their avatar’s movements. Eligible targets varied from one turn to the next but included all remaining vegetables on the farm (if the participant’s backpack was not full), the barn (if they were carrying one or more vegetables), or a small pillow icon labeled “(Pass)” to pass their turn (see top right of Fig. 1). Participants could not move to arbitrary open parts of the farm. A text box on the top left of the screen provided information to both players about whose turn it was; if a participant chose to pass their turn, or did not decide within 10 seconds, the text box briefly (2500ms) displayed why the turn had ended.

After an eligible farm target was selected, the player’s avatar automatically walked to the selected destination using the shortest path available to standardize action costs across participants (i.e., agents did not have to plan or optimize their walking path). The shortest path algorithm navigated around the location of other players because agents could not occupy the same grid tile at the same time. When participants visited the barn, all the objects they were currently carrying were deposited, leaving their bag empty. In addition, players immediately moved two tiles out of the way of the door to the barn so as to not block the entrance for the other player (these automatic steps incurred standard energy walking costs but were equal for each agent).

A score board was present on the right hand side of the screen. The amount of information shown varied by condition, but both players could always see the size and contents of the backpacks and the farm box (barn), their own remaining energy units, and their current score (a count of how many vegetables of their color were stored in the barn). The experiment manipulated whether players saw their partner’s current energy and score in addition to their own.

Experimental Conditions and Hypotheses

In order to provide a broad survey of the potential questions this paradigm can address, we varied several features of the environment and gameplay within and between subjects. First, we expected that aspects of the layout and design of the environment might influence helping decisions. Certain arrangements of objects in the environment might have

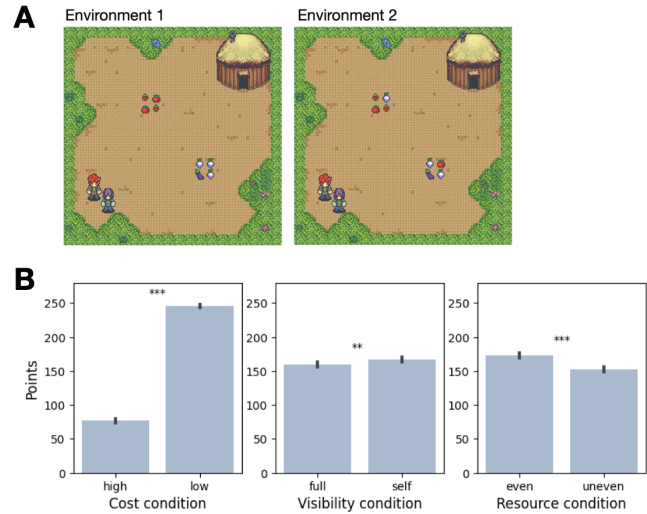


Figure 2: A: Two of the twelve starting arrangements that vary in “patch uniformity,” or whether adjacent vegetables are (Env. 1) or are not (Env. 2) of the same color. B: Performance by condition; average points earned in a game by a single subject. Error bars are 95% confidence intervals. Cost, visibility, and resource condition were constant within dyads across all twelve games in a session. (**: $p < 0.01$, ***: $p < 0.001$ from independent two-sample t -tests between conditions.)

made it easier, and more tempting, to help the other agent due to changes in the moment-to-moment energy cost of helping. Environment layout was varied within-subject such that participants encountered twelve unique environment layouts during the task.

In addition to the within-subject environment variations, we manipulated three two-dimensional task variables *between* dyads (cost, resource balance, and score visibility). We anticipated that the task’s cost structure, and specifically the overall global cost of helping, would impact collaborative behaviors (with helping being more common in cases where it carried lower personal cost). Next, we thought the balance of resources or ability amongst the two players would change their dynamics (as a person helping another may have some unique capacity lacked by the person needing help). Lastly, we considered that participants might help more or less depending on how much information was available to them about the score of their partner; participants might help more if the impact on their partner’s score was more salient.

The full factorial design incorporating these three two-dimensional variables therefore included $2^3 = 8$ between-dyad conditions. The following subsections detail each of the factors manipulated in the study.

Environment layout Of the twelve games each dyad played together, each game started with one of twelve unique “environments” reflecting the initial arrangement of farm items on the field. Players’ starting positions and the location

of the farm box were constant across all environments, but vegetable amounts and locations differed. Some example environments are shown in Fig. 2A, and all twelve environments are viewable on the project website. Environments were presented in a random order for each session. There were several features of the initial environment set-ups designed to elicit variation in behavior. In some games there were an equal number of the different colored vegetables, and in others there were more vegetables of one color. Sometimes the vegetables were located in patches of a single color, and other times both colors were mixed within a single cluster, which we called “patch uniformity” (as in Fig. 2A). The vegetable clusters could have different sizes and locations relative to the players. The range of environments was selected intuitively to test informal predictions that increasing the number of the partner’s vegetables and decreasing patch uniformity would increase helpfulness (because a player will more often be near their partner’s vegetable, allowing for less costly opportunities to help).

Cost - Low or High Players began with 100 energy points and used energy to walk around the grid world, with longer distances requiring more energy. Since the bonus point calculation for a given player was a product of the number of their vegetables harvested and their remaining energy, energy costs directly impacted the bonus payment participants could earn. The cost for walking was 1 energy unit per grid tile in the Low cost condition versus 2 energy units per tile in the High cost condition. In both conditions, passing one’s turn or not responding within ten seconds cost 5 energy units (no character movement). The energy level could not go beneath 0, and when energy was 0, the participant could still move around but would earn no bonus points on that game round.

Resource capacity - Even or Uneven The resource condition determined whether the two players had equal or unequal backpack sizes. In the Even resource condition, both participants had a backpack that could maximally hold 4 vegetables at a time. In the Uneven resource condition, one participant had a larger backpack (5 vegetable capacity) and one participant had a smaller backpack (3 vegetable capacity). The assignment of larger or smaller backpacks for dyads in the Uneven condition was random, and participants maintained the same backpack size across all twelve games in a session.

Energy visibility - Full or Self The visibility condition manipulated how much information about scoring and energy of the other player was available to each participant within the game display. In the Full visibility condition, participants saw the current energy level and score of both players in the game. In the Self condition, participants saw only their own energy and score. At the end of each game, the display showed either the bonus points earned by both participants (Full), or only the bonus points earned by the participant and not what their partner earned (Self).

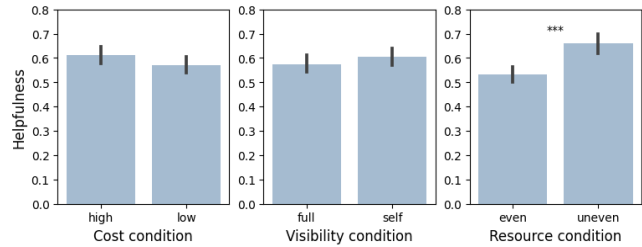


Figure 3: Helpfulness, or average number of helping actions in a game by a player, by experimental condition (as in Fig. 2B).

Statistical analyses

In addition to pre-registered tests of our experimental manipulations, we also report the results of several mixed regression models that helped account for nuance in the final data-set. Specifically, the pre-registration did not anticipate that the dependent variable was not normally distributed. As such, we turned to models with less strict assumptions. Generalized linear mixed models were fit in R with the package *lme4* (Bates, Mächler, Bolker, & Walker, 2014) and were constructed by starting from a minimal model, then adding task-relevant parameters only when they improved the model fit according to likelihood ratio tests.

One mixed effects linear regression model, which we label the Game-level Performance Model, predicted subjects’ *performance* (i.e., total earned points) in a *game* given various contextual features. The model consisted of eighteen fixed effects parameters which accounted for environmental features (e.g., number of vegetables in the initial layout) and recent helpful behavior (e.g., whether the player’s partner helped in the previous game).

Separate but similar models predicted the likelihood for a participant to *help* their partner at a *game* level (i.e., did the participant help at any point in this game round?) and at a *trial* level (i.e., did the participant help on this specific trial?). These models had sixteen parameters since helpfulness was the variable being predicted rather than a variable affecting performance. The Game-level Helping Model provided a measure of global features that encouraged or discouraged helping, while the Trial-level Helping Model identified local contextual factors that led to a specific helpful decision. All models included random intercepts to account for differences between subjects.

Results

After exclusions due to incomplete data, our dataset consisted of $N = 628$ participants (314 dyads). The sample was predominantly male-identifying (372 male, 243 female, 13 other/no response), and 76% of participants indicated their race as Caucasian/white. The average age of participants was 37 years old ($SD = 11.7$). The game data are rich consisting of moment-by-moment choices of individuals in a dyadic task, but for the purposes of this analysis we focus on exam-

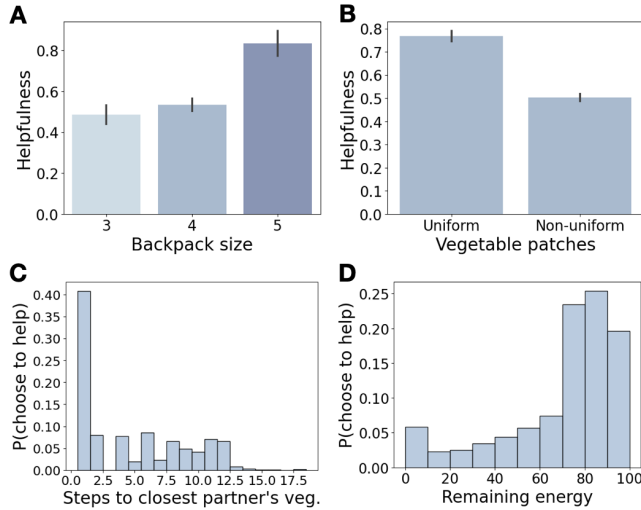


Figure 4: A and B show how “helpfulness” (average number of helping actions in a game by a player) varies with player backpack size and the color uniformity of clusters (patches) of vegetables, respectively. C and D present the proportion of trials in which a player chose to help given local contextual features; In C, the costliness of the least costly helpful action measured as distance (number of steps) to the closest partner’s vegetable, and in D, the player’s number of remaining energy units which contributed to bonus calculation.

ining how well participants performed the task and how much “help” they provided in different contexts.

Pre-registered predictions

We conducted pre-registered Bayesian *t*-tests to evaluate the effect of the between-dyad manipulations. Our data provided support for only one of our predictions; a one-sided *t*-test showed participants helped more when they had a larger capacity backpack than their partner ($t(288) = 3.07, p < .01, BF10 = 22.2$). However, the dependent variable was not normally distributed, so these tests were not conclusive.

Game-level Performance

Fig. 2B shows average participant performance in several of the experimentally manipulated conditions. While several qualitative trends are readily apparent at first glance, the regression model provided more rigorous testing to back up our conclusions. We report regression parameter estimates, along with 95% confidence intervals (*CI*) and *p*-values ($Pr(> F)$). While space prohibits a full presentation of the model results, one key finding was that participant performance decreased in games where they provided help ($-5.91, CI = [-9.33, -2.47], p < .001$) and increased with help from their partner in the previous game ($5.69, CI = [2.35, 9.06], p < .001$) and current game ($30.48, CI = [27.06, 33.94], p < .001$). Critically, this result confirms that helping in the Farm Task is costly but does benefit the player being helped.

Game-level Helping

Most importantly, we were interested in quantifying when participants engaged in helping as a function of the features of the environment. Our pre-registration defined helping as the amount of energy a player expended to pick up vegetables of their partner’s color. Although each player could pick up and deliver vegetables of both colors, their reward depended only on their own color items and energy. Therefore, when a player chose to pick up a vegetable of their partner’s color, they sacrificed their own energy without direct compensation.

The Game-level Helping model predicted the likelihood that a player helps in a game given contextual features. Of particular interest were features that reflected the amount of helping of one’s partner in the current and previous game, as well as one’s own helpfulness in the previous game. We report parameter estimates β , their 95% *CI*s, and computed Odds Ratios (*OR*) to aid interpretation of how each feature increased ($OR > 1$) or decreased ($OR < 1$) the likelihood of a participant helping in a game.

Although we had initially predicted that helping would occur more often as action costs decreased, the fitted model showed no effect of the global cost condition on helping in this experiment design. We cannot conclude from this that costs do not matter to helping, but only that this design was insufficient to detect a between-dyad effect. In contrast to our predictions that helping would increase when one’s partner’s score was visible, and that helping would increase when resources (backpack capacity) were unevenly distributed, the results showed no effect of the visibility ($t(305) = -0.22, p = .59, BF10 = 0.26$) or resource ($t(306) = 0.61, p = .27, BF10 = 0.30$) conditions on helping at the between-dyad level at which these conditions were varied (see Fig. 3 for game-level data). However, we did find that the likelihood of helping increased as a player’s backpack size increased ($OR = 3.82, \beta = 1.34, CI = [1.12, 1.56]$) (also see Fig. 4A). That is, in the uneven resource condition, the player with the larger backpack was more likely to help than their partner.

The initial arrangement of each game affected how much participants helped. When there were more of a player’s own vegetables in the starting environment, they helped less often ($OR = 0.22, \beta = -1.50, CI = [-1.61, -1.38]$), whereas they helped more often when there were more of their partner’s vegetables ($OR = 4.66, \beta = 1.54, CI = [1.41, 1.67]$). Although we might trivially expect that players pick up their partner’s vegetables more often when more of their partner’s vegetables are present, players could also choose to pass their turn instead of helping. Instead, players more often helped their partner, perhaps in an effort to finish the round more quickly than if they were inactive. As shown in Fig. 4B and bolstered by the modeling results, participants helped less when the patches were all uniform ($OR = 0.25, \beta = -1.37, CI = [-1.59, -1.15]$), indicating a division of labor segregated by clusters rather than color. That is, players prioritized minimizing the distance between their actions over picking up only their own color vegetable on each action.

Trial-level Helping

The Trial-level Helping model predicted the likelihood of a subject choosing a helpful action on their turn. For a trial-level model, there are a large number of potentially relevant parameters, since it is possible that all aspects of the current game state, such as the location of every vegetable or the players' score and energy status, ultimately affect the decision of whether to help. For our model, we focused on features that would be relevant in the construction of a cognitive model, including the experimental manipulations, the player's remaining energy, local costs of helping, and recent helpfulness. We expected that players would help more often when a helpful action was less costly, and that they would also help more often if they considered their partner to be "helpful."

The local cost of helping was defined as the costliness of the cheapest helpful action a player could take on that trial - specifically, the Manhattan distance from the player to their partner's closest vegetable (Fig. 4C). As we expected, the likelihood of helping decreased as the local cost of helping increased ($OR = 0.45$, $\beta = -0.80$, $CI = [-0.84, -0.75]$). When the cheapest helpful action available was more costly, participants helped less often.

Players took into account their own capacity when deciding whether to help. Specifically, players were more likely to help when they had more "energy" units remaining ($OR = 1.27$, $\beta = 0.24$, $CI = [0.13, 0.35]$) (Fig. 4D). Taken together with the positive effect of backpack size on helping (Fig. 4A), these results indicate that participants are more willing to help when they have greater capacity, and might not help if resources seem scarce.

Measures of reciprocity had a large impact on whether a player decided to help. Players were more likely to help when their partner had helped on the previous trial ($OR = 1.35$, $\beta = 0.30$, $CI = [0.27, 0.33]$) and as the total number of helping events by their partner across the whole experiment increased ($OR = 1.26$, $\beta = 0.23$, $CI = [0.17, 0.30]$).

Comparison with a heuristic agent

To further examine trial-level helping and direct sequential reciprocation (e.g., one turn to the next), we simulated behavior from an agent following a simple decision-making heuristic. The agent followed a Nearest Neighbor policy, always selecting the closest vegetable. Fig. 5 shows the averaged proportion of helpful actions on each turn within a game in the human data (left) and in simulated data between two Nearest Neighbor agents (right). Separate curves distinguish whether or not the deciding player's partner had helped on the previous turn. The human data, unlike the simulated data, shows a large gap between the two curves, revealing how the likelihood to help changes drastically as a function of whether one's partner has recently helped (potentially ruling out this is an artifact of other aspects of the game).

Discussion

In this work, we conducted a large, pre-registered, factorial experiment designed to explore human helping behavior. An

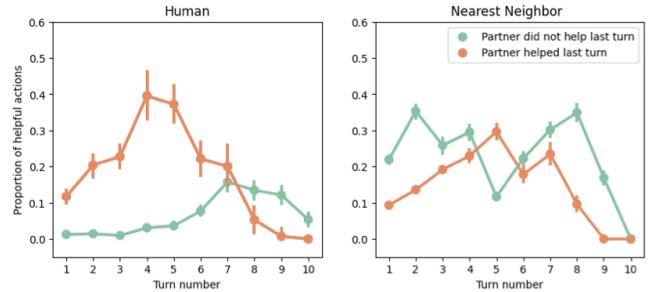


Figure 5: Proportion of helping events on each trial (turn) averaged across subjects. The orange curve reflects the proportion of helping actions selected on trials where the player's partner had helped on the previous trial, while the teal curve reflects helping when the partner did not help on the previous trial. The proportion is normalized by the number of helpful actions available on each turn. Plots reflect data from the first ten turns of games beginning with eight vegetables.

online, interactive two-player game was designed to examine when, why, and how people help others. One amazing aspect of our design is that helping other participants was not explicitly incentive-compatible in the design (people were only paid for vegetables of their own color). As a result, all the instances of helping in the task reflect genuine examples of players spontaneously deciding to help their partner.

The results showed three key findings. First, people were more likely to help when they were endowed with greater resources than another agent. This effect was larger than the manipulations of overall cost and score visibility. Second, helping was inherently reciprocal, as participants help another when they come to expect help for themselves from the other player. However, people also think about their own opportunity costs and benefits and are unlikely to help if the cost of helping is large, as evidenced by the large effect of the layout of the task environment on the moment-to-moment decision to help.

While this work reveals local contextual features that affect participants' momentary decisions to help, the Farm Task design allows for many further investigations into helping and collaboration. Behavior was non-trivial and could not be explained by a simple self-serving heuristic model, indicating that the environment allows for complex cognition and planning. Future work can construct cognitive models using the relevant task variables to build generative models of behavior in similar environments. Comparing human behavior to generative agents can also refute alternate explanations for helping behavior, such as desire to finish the task quickly.

The virtual environment provided by the Farm Task allows us to identify subtleties of helping behavior, highlighting decisions that fall somewhere in between purely selfish and altruistic. Future work aims to build upon these results to identify computational models of how people decide when and how to help.

Acknowledgments

The authors thank the National Science Foundation for funding through the project "Towards a computational cognitive science of helping," Award No. 2021060.

References

- Andreoni, J., & Miller, J. H. (2008). Analyzing choice with revealed preference: is altruism rational? *Handbook of experimental economics results, 1*, 481–487.
- Bates, D., Mächler, M., Bolker, B., & Walker, S. (2014). Fitting linear mixed-effects models using lme4. *arXiv preprint arXiv:1406.5823*.
- Boyd, R. (1988). Is the repeated prisoner's dilemma a good model of reciprocal altruism? *Ethology and Sociobiology*, 9(2-4), 211–222.
- Bridgers, S., & Gweon, H. (2018). Means-inference as a source of variability in early helping. *Frontiers in Psychology*, 9, 1735. doi: 10.3389/fpsyg.2018.01735
- Fehr, E., & Fischbacher, U. (2003). The nature of human altruism. *Nature*, 425(6960), 785–791.
- Fehr, E., & Gächter, S. (2000). Fairness and retaliation: The economics of reciprocity. *Journal of economic perspectives*, 14(3), 159–182.
- Ho, M. K., MacGlashan, J., Littman, M. L., & Cushman, F. (2017). Social is special: A normative framework for teaching with and learning from evaluative feedback. *Cognition*, 167, 91–106.
- 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 *Cogsci*.
- Marsh, A. A. (2016). Neural, cognitive, and evolutionary foundations of human altruism. *Wiley Interdisciplinary Reviews: Cognitive Science*, 7(1), 59–71.
- Palan, S., & Schitter, C. (2018). Prolific.ac—a subject pool for online experiments. *Journal of Behavioral and Experimental Finance*, 17, 22–27.
- Stephens, C. (1996). Modelling reciprocal altruism. *the British Journal for the Philosophy of Science*, 47(4), 533–551.
- Tomasello, M. (2009). *Why we cooperate*. MIT press.
- Trivers, R. L. (1971). The evolution of reciprocal altruism. *The Quarterly review of biology*, 46(1), 35–57.
- Von Ahn, L., Blum, M., Hopper, N. J., & Langford, J. (2003). Captcha: Using hard ai problems for security. In *Eurocrypt* (Vol. 2656, pp. 294–311).