

Two's company, six is chaos: conventions emerge regardless

Anonymous CogSci submission

Abstract

TODO In repeated reference games where a speaker describes the same set of images to a listener over a series of rounds, the number of words used decreases as the pair converge on ad-hoc names for the images. The dynamics of this efficient reference formation is well-studied in dyads; however much communication takes place in larger groups, which are rarely studied in this paradigm. The current work extends iterated reference games to groups of 2-6 people who rotate between speaker and listener roles in an online game with text-based communication. Across TODO games and more than TODO K total words, we find high accuracy and patterns of reduction regardless of group size.

Keywords: TODO; Add your choice of indexing terms or keywords; kindly use a semi-colon; between each term.

TODOs

Do we include fun anecdotes?

check what current empirica citaiton is

facts to include: num games, num words overall

Intro

TODO re-write intro! [blah blah communication in non-dyads] Verbal communication is an integral part of our daily lives. We coordinate schedules with partners, socialize with friends over board games, learn and teach in seminar classes, and listen to podcasts. Our communicative environments range in size from one-on-one to small group to large group to broadcast, but the goal of efficient communication is held in common. One necessity for efficient communication is shared reference expressions; when referring to a thing or an idea, it needs some sort of name that the interlocutors will jointly understand. In many cases, there are widely shared conventionalized expressions, but in other cases, spontaneous ad-hoc expressions are needed.

[hole to fill] The formation of these new reference expressions is well-studied in dyadic contexts; however, dynamics may be different in larger groups, which are less studied. Our current work builds on the dyadic reference game tradition by extending it to small groups.

Dyadic reference games

The typical paradigm for studying partner-specific referring expressions is an iterated reference game with asymmetric knowledge. That is, each round there is a speaker who knows

what the target is and a listener who does not have this information. In Clark & Wilkes-Gibbs (1986), each speaker described 12 tangrams in order, so their listener could correctly order the images. After receiving feedback, the pair repeated the task with the same images but a new order, for a total of 6 repeats. Crucially, Clark & Wilkes-Gibbs (1986) found that reference expressions condense over the course of repeated reference to the same image. Early descriptions are long and make reference to multiple features in the image, but by later rounds definite shorthand names are used.

Recently, online participant recruitment and web-based experiments have made it possible to study this convergence in larger populations using a text-based communication interface. In Hawkins, Frank, & Goodman (2020), 83 pairs completed a cued version of the iterated reference experiment. On each trial, the speaker saw one image highlighted and described it to the listener who clicked on what they thought the target was. Both players received feedback before moving on to the next target image. All images were highlighted each block, for a total of 6 blocks. Speakers produced fewer words per image in later blocks than in earlier blocks, in line with results from face-to-face, oral paradigms.

While this reduction pattern is robust for dyads, less is known about the how utterances are adapted in larger groups. A couple of studies point to some potential difficulties in trying to communicate with multiple people at once.

Multi-party communication

Yoon & Brown-Schmidt (2019) had speakers complete a sorting task with some listeners, so that they would have a common ground of shared names for the images. Then in a test phase, the speaker described these images to a group of either all knowledgeable listeners from the sorting task, new listeners who had not done the sorting task, or a mix of knowledgeable and new listeners. Speakers produced longer utterances when any new listeners were present than with only experienced listeners. This might predict slower reduction in larger groups where there will inevitably be some variability in how people understand reference expressions. These studies included 3-hour experiments that were very time and labor intensive, but some of the questions about group dynamics may be addressable in online experiments taking advantage of natural variation in understanding without artificially inducing large knowledge differences.

It’s difficult to communicate with naive listeners, but it can be even harder to communicate with someone with entrenched preconceptions. Weber & Camerer (2003) induced these conceptual differences by having two pairs of people (each pair representing a “firm”) do an iterated reference game with the same set of pictures. After 20 rounds, there was a “merger” where the listener from one group joined the other group. The reference game continued with the speaker communicating to both their original listener and the new listener. After the merger, there was a jump in how long it took either listener to make a selection. Even after several more rounds, listeners were still not as fast as before the merger. With larger groups of people all speaking together, there’s a greater chance for different people to independently develop different conceptualizations of an image, and it may be difficult for them to understand each other or agree on a common term of reference.

But what’s the dose-response curve?

Conventions across speakers

[set-up that there’s an expectation of consistency w/i speaker except when new people join – what is it when everyone’s been together?] Studies vary in whether the same participant keeps the speaker role the entire game (ex. Clark & Wilkes-Gibbs, 1986) or whether the roles rotate (ex. pre-merger rounds of Weber & Camerer, 2003). Role rotation makes the paradigm more similar to collaborative puzzle-solving exercises also used to study conventions (Garrod & Doherty, 1994; Ibarra & Tanenhaus, 2016).

In general, listeners expect conventions to be maintained, but they are not surprised to new descriptions of a familiar object if it comes from a new speaker (Metzing & Brennan, 2003) or if a new listener is present (Yoon & Brown-Schmidt, 2014). It’s unclear how these expectations map onto a group of people rotating roles in the task who are all present the entire time. Do later speakers count as new, or are they expected to follow conventions they’ve already heard? Do additional non-new listeners license changes in descriptions?

Present work

In this work, we extend the dyadic repeated reference game paradigm of Hawkins et al. (2020) to games for 2-6 players who rotate between speaker and listener roles. We compare accuracy and reduction rates in groups of different sizes.

Methods

TODO methods summary! We adapted the methods of Hawkins et al. (2020), adjusting them to work for multi-player games with rotating speakers. Participants played a repeated reference game where a speaker saw an array of tangrams with one indicated (Fig 1 and had to communicate which figure to click to the listeners using the chat box. Within each block, each of the 12 tangrams was the target once, and the speaker role rotated each block, so all participants were the speaker at least once. Games ran for a total of 6 blocks. We recorded what participants said in the chat,

as well as who selected what image and how long they took to make their selections. The experiment was implemented in Empirica (Almaatouq et al., 2020); materials to run the experiment, as well as data and code are available at TODO anonymous OSF clone .

We report how we determined our sample size, all data exclusions (if any), all manipulations, and all measures in the study. Our preregistration is at <https://osf.io/cn9f4>. TODO check anonymity

Participants

Players	Partial	Complete
2	4	15
3	2	18
4	2	19
5	3	17
6	6	12

Table 1: Number of games run for each player count.

Participants were recruited using the Prolific platform between May and July 2021. We screened for participants who were fluent, native English speakers. Participants were paid \$7 for 2-player games, \$8.50 for 3-player games, \$10 for 4-player games, and \$11 for 5- and 6-player games (with the intention of a \$10 hourly rate), in addition to up to \$2.88 in performance bonuses.

Materials

We used the 12 tangram images used by Hawkins et al. (2020) and Clark & Wilkes-Gibbs (1986) (see Fig 1). These images were displayed in a grid with order randomized for each participant. The same images were used every block.

Procedure

We implemented the experiment using Empirica, a Javascript-based platform for running real-time interactive experiments online (Almaatouq et al., 2020). Code for running this experiment is available at TODO anonymize <https://github.com/vboyce/FYP>. From Prolific, participants were directed to our website where they navigated through a self-paced series of instruction pages explaining the game. Participants had to pass a quiz to be able to play the game. They were then directed to a “waiting room” screen until their partners were ready.

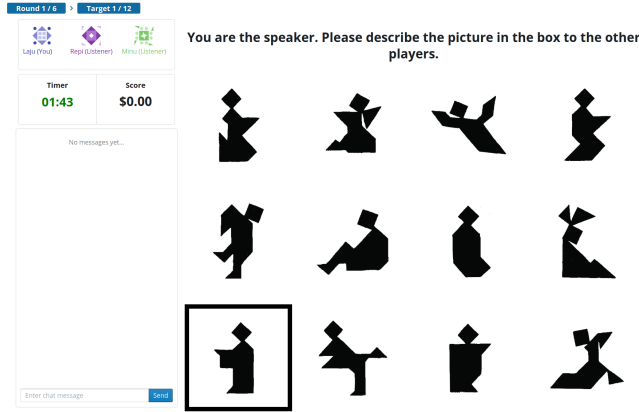


Figure 1: Screenshot of the speaker's view. Participants see all 12 tangram images.

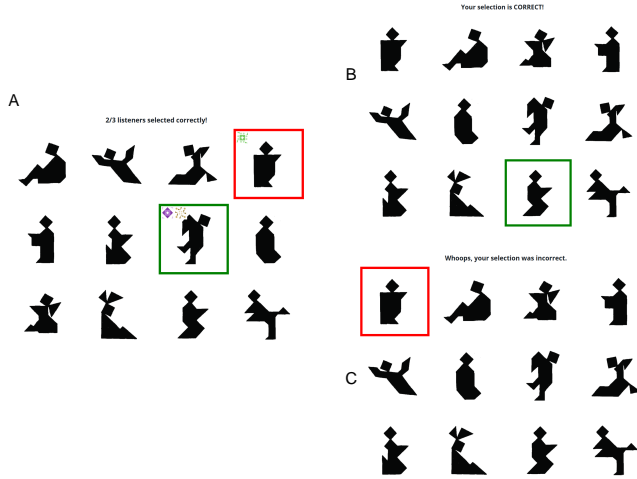


Figure 2: Screenshots of feedback for speakers and listeners. Speakers (A) saw what figure each person chose, indicated by the matching icons. Listeners only learned if their selection was correct (B) or incorrect (C). Listeners were not shown what other listeners chose.

Once the game started, participants saw screens like Fig 1. Each trial, the speaker described the highlighted tangram image so that the listeners could identify and click it. All participants were free to use the chat box to communicate, but listeners could only click once the speaker had sent a message. Once a listener clicked, they could not change their selection. Once all listeners had selected (or a 3-minute timer ran out), participants were given feedback. Listeners only learned whether they individually had chosen correctly or not; listeners who were incorrect were not told the correct answer. The speaker saw which tangram each listener had selected. Listeners got 4 points for each correct answer; the speaker got points equal to the average of the listeners' points. These points translated into cents of performance bonus at the end of the experiment.

In each block, each of the 12 tangrams was indicated to the speaker once. The same person was the speaker for an entire block, but participants rotated roles between blocks. Thus, over the course of the 6 blocks, participants were speakers 3 times in 2-player games, twice in 3-player games, once or twice in 4 and 5-player games, and once in 6-player games.

After the game finished, participants were given a survey asking for optional demographic information and feedback on their experience with the game.

Data pre-processing

Participants could use the chat box freely, which meant that the chat transcript contained some non-referential language. The first author skimmed through the chat transcripts, tagging utterances that did not refer to the current tangram. These were primarily pleasantries ("Hello"), meta-commentary about how well or fast the task was going, and confirmations or denials ("ok", "got it", "yes", "no"). We exclude these utterances from our analyses. Note that chat lines sometimes included non-referential words in addition to words referring to the tangrams ("ok, so it looks like a zombie", "yes, the one with legs"); these lines were retained intact.

Our intended sample size was 20 complete games in each group size, but we ended up with fewer as shown in Table 1. We excluded incomplete blocks from analyses, but included complete blocks from partial games. (Partial games occurred when a participant disconnected early, for example due to internet trouble.)

Results

Across groups of all sizes, participants did well at the task, accurately identifying the tangrams most of the time. In later blocks, participants were faster, and fewer words were used to achieve increasing levels of accuracy. This replicates the classic pattern of reduction found in dyadic versions of the game.

Speed and Accuracy

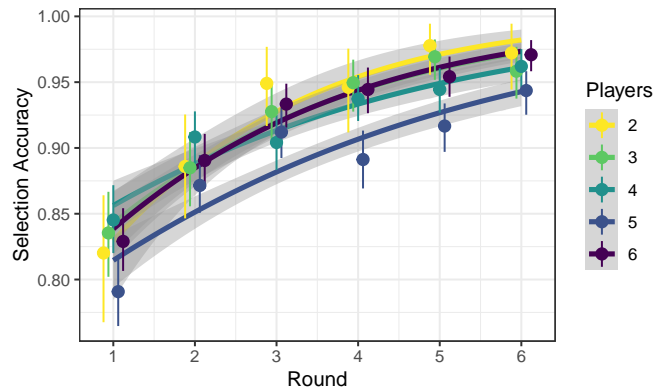


Figure 3: Accuracy TODO

Most groups were accurate in their selections, with accuracy rising over blocks (Fig 3). Participants are more accurate in later blocks (block: $\text{Est}=0.38$, $\text{CrI}=[0.25, 0.5]$), and there was no strong effect of group size on accuracy (numPlayers: $\text{Est}=-0.02$, $\text{CrI}=[-0.08, 0.03]$).

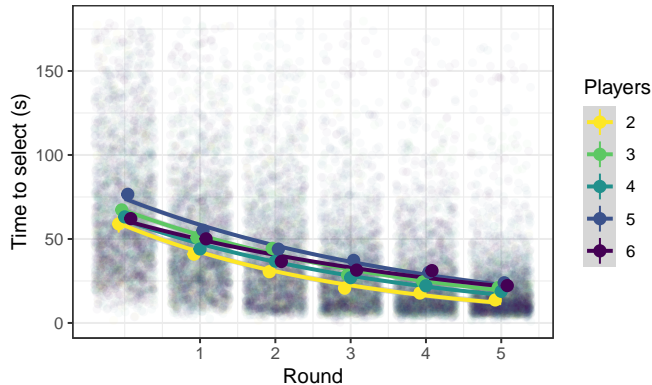


Figure 4: Listeners selected images faster in later blocks. Only times to correct responses are shown.

Participants selected images faster in later blocks (Fig 4), although there is wide variability. In a linear model, the participants are faster each block (block: $\text{Est}=-10.03$, $\text{CrI}=[-11.03, -9.03]$) and are slightly slower in larger games (numPlayers: $\text{Est}=1.03$, $\text{CrI}=[0.4, 1.66]$). This speed up is consistent with prior work by Weber & Camerer (2003) which used speed as the dependent measure. [? TODO stats]

Convention

Both speakers and listeners reduce the amount they say over the course of blocks. As shown in Fig 5, the number of words produced does decrease. Listeners often don't talk much, but are more likely to ask questions or make clarification in early blocks. In a linear regression for the number of words each listener said, they say less in later blocks (block: $\text{Est}=-0.48$, $\text{CrI}=[-0.79, -0.18]$), but there is no clear effect of game size (numPlayers: $\text{Est}=0.2$, $\text{CrI}=[-0.13, 0.51]$).

Speakers make longer utterances in early blocks and reduce to shorter utterances in later blocks. In aggregate, the effect of being one block later is block: $\text{Est}=-3.35$, $\text{CrI}=[-4.58, -2.13]$. The overall effect of having more players in a group is numPlayers: $\text{Est}=1.67$, $\text{CrI}=[0.68, 2.71]$ per additional player. Notably, this shortening pattern occurs even as speakers rotate. This estimate is uncertain because of a relatively small number of groups and wide group-level variability. TODO modelling!

Because the ground truth answers are not provided to listeners who make mistakes, they may not learn what an utterance referred to (unless they ask in the chat). What happens if a listener gets a tangram wrong and then is the speaker on the next block? For that tangram, they are unlikely to build off the previous description they didn't understand. In contrast, a speaker who previously got the tangram right is likely

to continue the conceptualization used so far and conventionalize it more, such as by reducing unneeded details. Taken together, this leads to the hypothesis that speakers should say more words when they got the tangram wrong the previous block, after controlling for other effects. This is borne out; speakers say was_INcorrect: $\text{Est}=3.07$, $\text{CrI}=[1.65, 4.5]$.

Variability

While the averaged data shows a smooth reduction in the number of words, individual trajectories for specific tangrams in specific groups are more varied. Reduction is not monotonic, as some later speakers use more words than were used in earlier blocks.

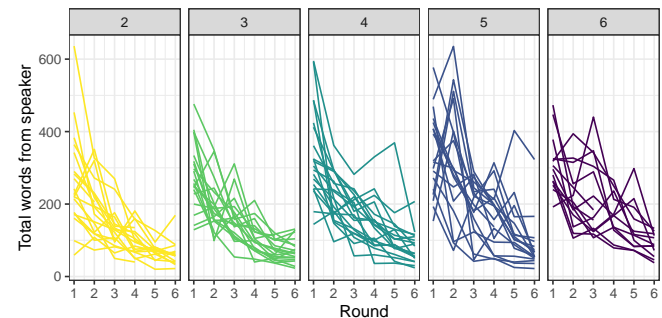


Figure 6: Number of words from speaker (total, across all 12 figures) in a block. Convention formation is not a smooth process, and some later speakers in some groups used more words than their earlier speakers.

Content similarity

TODO might want to pick better analyses??

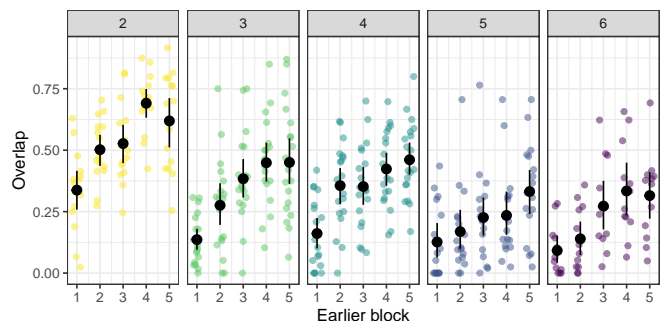


Figure 7: The fraction of content words used by the speaker in the last block that were used to describe the same tangram in an earlier block. TODO better caption

TODO the brm won't run so using glm, but that means we cache and also I need to figure out how to extract numbers from it.

Discussion

The overall pattern of utterances shortening over repeated reference extends to groups of 3 or 4 people talking together and

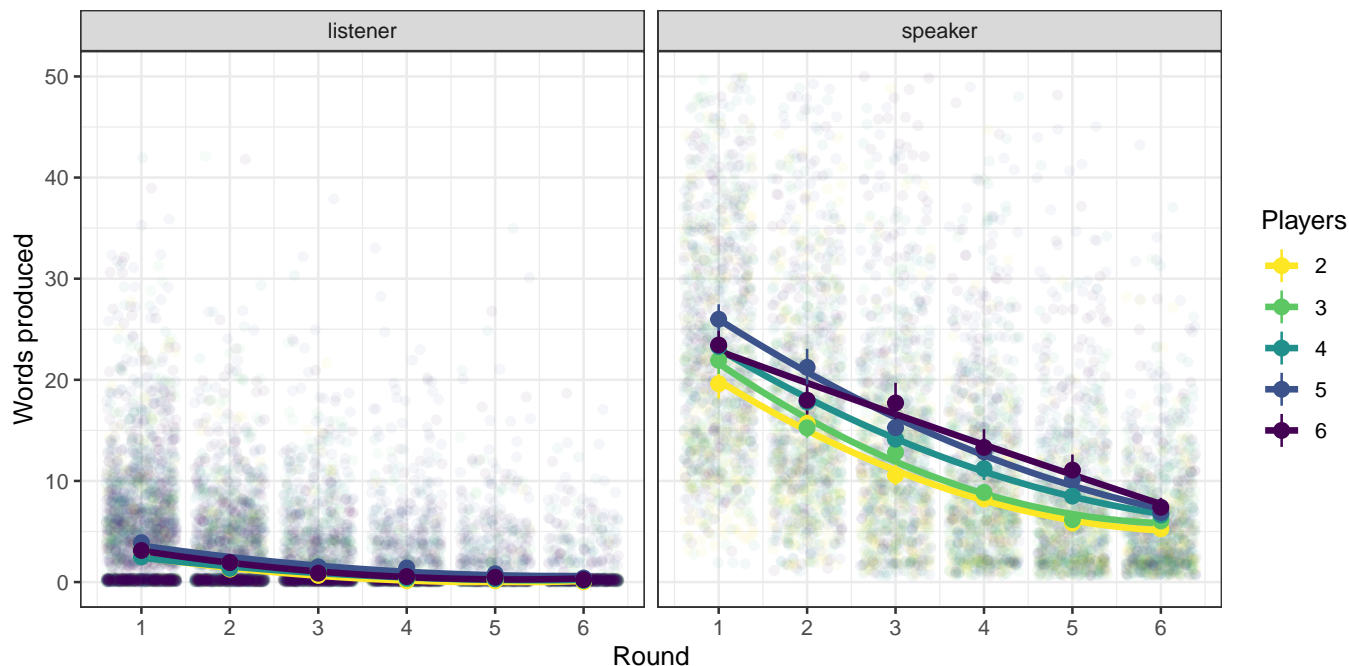


Figure 5: Speaker and listeners say fewer words in later blocks. Note: y-axis clipped at 50 which hides a few speaker outliers.

rotating between speaker and listener roles. Rotating speakers gives a stronger interpretation of reduction as conceptual agreement because more people have to produce the short-hand names.

We provided less feedback than previous studies such as Hawkins et al. (2020). This low level of feedback means that there isn't a way for people to find out what was meant for utterances they initially did not understand outside of the verbal communication channel (or process of elimination). Similarly, speakers don't have direct access to how well their partners did in the previous block. Real-life communicative situations vary in what extra-textual feedback exists, but we do show that people can work around their initial confusion to eventually understand utterances, rather than just memorizing pairings after the first occurrence.

This is a rich data set consisting of TODO words across TODO referring expressions by TODO speakers, in addition to clarifications questions and comments from listeners. In this set of analyses, we rely on the easy to calculate measures of accuracy, speed, and word counts as proxies for the content of the utterances. In future analyses, it would be useful to do content analysis to understand how and when concepts are introduced and conventionalized and how much the semantics of utterances varies block to block (and speaker to speaker) depending on group size.

We demonstrate that it is feasible to extend iterated reference game paradigms to small groups of participants using an online platform, and thus rapidly gather high quality utterance data from a number of games. We found that the widely observed pattern of partner specific adaptation and reduction extends to 3 and 4 person games. Inter-group vari-

ability suggests that a closer look at interpersonal communication dynamics, for example, comparing the semantic content of utterances of players in the same or different games is warranted. A closer analysis of the utterances may yield information about how humans adapt language quickly, and the dataset may be useful for training artificial agents to use and understand language more dynamically.

justify my low info high rotation format # Acknowledgements

Place acknowledgments (including funding information) in a section at the end of the paper.

References

- Almaatouq, A., Becker, J., Houghton, J. P., Paton, N., Watts, D. J., & Whiting, M. E. (2020). Empirica: A virtual lab for high-throughput macro-level experiments. *arXiv:2006.11398 [Cs]*. Retrieved from <http://arxiv.org/abs/2006.11398>
- Clark, H. H., & Wilkes-Gibbs, D. (1986). Referring as a collaborative process. *Cognition*.
- Garrod, S., & Doherty, A. (1994). Conversation, coordination and convention: An empirical investigation of how groups establish linguistic conventions, 12.
- Hawkins, R. D., Frank, M. C., & Goodman, N. D. (2020). Characterizing the dynamics of learning in repeated reference games. *arXiv:1912.07199 [Cs]*. Retrieved from <http://arxiv.org/abs/1912.07199>
- Ibarra, A., & Tenenhaus, M. K. (2016). The Flexibility of Conceptual Pacts: Referring Expressions Dynamically Shift to Accommodate New Conceptualizations. *Front. Psychol.*, 7. <http://doi.org/10.3389/fpsyg>

.2016.00561

- Metzing, C., & Brennan, S. E. (2003). When conceptual pacts are broken: Partner-specific effects on the comprehension of referring expressions. *Journal of Memory and Language*, 49(2), 201–213. [http://doi.org/10.1016/S0749-596X\(03\)00028-7](http://doi.org/10.1016/S0749-596X(03)00028-7)
- Weber, R. A., & Camerer, C. F. (2003). Cultural Conflict and Merger Failure: An Experimental Approach. *Management Science*, 49(4), 16.
- Yoon, S. O., & Brown-Schmidt, S. (2014). Adjusting conceptual pacts in three-party conversation. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 40(4), 919–937. <http://doi.org/10.1037/a0036161>
- Yoon, S. O., & Brown-Schmidt, S. (2019). Audience Design in Multiparty Conversation. *Cognitive Science*, 43(8), e12774. <http://doi.org/10.1111/cogs.12774>