

Interaction structure constrains the emergence of conventions in group communication

Veronica Boyce^{1,*}, Robert Hawkins², Noah D. Goodman¹, Michael C. Frank¹

¹Stanford University

²Princeton University

Abstract

Multi-party communication is ubiquitous, but the need to take multiple interlocuter’s perspectives into account presents challenges not found in dyadic communication. One test case for communication is iterated reference games, where the phenomenon of reduction over repeated reference is well-attested in dyadic contexts and could explain how pairs of people build shared meanings. We extend the repeated reference game paradigm to groups of 2 to 6 people under varied interaction structure constraints across 313 games (1319 participants). Across conditions, groups shorten their utterances and form shared descriptions. Smaller groups and groups with thicker communication channels converge to shared conventions more rapidly than larger groups with thin communication channels. [TODO need another sentence translating this back to the unique challenges of multi-party communication and exposing the relevance of the results for theory/application (e.g. ‘Taken together, these results provide new insight into the conditions under which group communication can thrive, with implications for education, management, and ...’)]

TODO check ordering of legends on Figures!

TODO Figure 1 still under discussion!

intro could use more cites

cites to incorporate somewhere: adhoc conventions needed when priors are more different Murthy et al. (2022) —

While most psychological studies of communication focus on one-on-one conversations, much of our daily lives center around communicating in larger groups. In school, children sit in classrooms with peers and teachers; at home, we have dinner with multiple friends and family; and in the office, we attend meetings with colleagues and managers.

Communicating in groups can present challenges not present in two-party (dyadic) communication (Traum 2004, Ginzburg & Fernandez 2005, Branigan 2006). Speakers need to produce utterances for multiple listeners who may have different levels of background understanding (Horton & Gerrig 2002, Weber & Camerer 2003, Horton & Gerrig 2005, Fox Tree & Clark 2013, Yoon & Brown-Schmidt 2014, 2018). Listeners must comprehend utterances that are less tailored to them due to the presence of other interlocuters (Schober & Clark 1989, Carletta et al. 1998, Fay et al. 2000, Metzing & Brennan 2003, Rogers et al. 2013, Tolins & Fox Tree 2016, Cohn-Gordon et al. 2019, Yoon & Brown-Schmidt 2019). In spite of these considerable challenges, however, speakers are often able to navigate multi-party settings with ease. What aspects of communication channels and interaction structure determine how efficiently a group can communicate?

One component of effective communication is establishing shared reference, that is a joint understanding among interlocuters that a certain phrase uttered by the speaker refers to a certain object or idea in the world. Sometimes, shared reference is the end goal (in for example, “can you hand me the small bowl?”), and sometimes shared reference is a pre-requisite to larger conversational goals (example, in a debate, interlocuters still need a mutual understanding of what the words mean). In many cases, conventional semantic meanings of words from a shared language plus culturally shared conventions around language

*Corresponding author. Email: vboyce@stanford.edu

use are sufficient for two people to establish shared reference. In other cases, mappings between words and objects may rely on shared history to be understood (ex. “pass me my favorite cup.”).

Reference is a component of communication that can be isolated and tested in experimentally manipulated contexts, so it has been used as a case study for efficient communication more broadly. Some reference games examine how interlocutors communicate in the absence of conventional labels. In these cases, interlocutors must invent ad-hoc reference expressions to communicate about objects without canonical names.

The formation of these new reference expressions is well-studied in dyadic contexts. Clark & Wilkes-Gibbs (1986) established an experimental method for studying the emergence of new referring expressions that has now become standard (building on Krauss & Weinheimer 1964, 1966). Two participants see the same set of figures; the speaker describes each figure in turn so the listener can select the target from the set of figures. The speaker and listener repeat this process many times with the same images.

Early descriptions are long and make reference to multiple features in the figure, but over the course of the game, shorthand conventional names for each figure emerge; this shortening of utterances is called ‘reduction’. Not only are later utterances shorter than earlier utterances, but later utterances are a tacitly agreed upon name, understandable within the dyad, but different from the conventions chosen by other dyads.

Recently, online participant recruitment and web-based experiments have made it possible to study this convergence in larger populations (Haber et al. 2019, Hawkins et al. 2020). In line with results from face-to-face, oral paradigms, in these online, chat-based paradigms, speakers reduced their utterances, producing fewer words per image in later blocks than in earlier blocks. (Throughout this paper, we use “speaker” and “listener” to refer to the roles describing and selecting targets, regardless of communication modality.)

What aspects of conversational infrastructure are needed to support this reduction to efficient referring expressions? In the current work, we address how components of interaction structure, including group size and communication channels, shape how successfully groups form partner-specific conventionalized names for target objects over the course of an iterated reference game. We recruited 1319 participants who were organized into 313 groups distributed across 3 online experiments and 11 conditions. Collectively, players produced 326000 words during their games. We analysed the games along 4 metrics: two group performance metrics of 1) listener accuracy and 2) number of words produced per trial, as well as two computational measures of semantic similarity that addressed 3) how utterances converge towards a conventions with a game and 4) how utterances diverge from descriptions in other games as they become more partner-specific.

Results

Reference game experiments can vary along numerous dimensions: how many people are participating, whether they have any prior relationships with one another, what roles they have within the game, what types of targets they are describing, what type of communication channel they have, and what sorts of feedback they get on their selections, among others. This makes for a very large space of possible iterated reference game experiments (Almaatouq et al. 2022)!

We focused on dimensions that could affect how groups interacted with each other and sampled along four dimensions that parameterized the interaction space. Our experiments varied along the dimensions of game size, speaker rotation, level of feedback, and the form of the listener backchannel; other possible dimensions, such as what the target images were or if participants had prior relationships, were held constant across games. Across the three experiments, we sampled 11 different points in the experimental space, illustrated in Figure 1A.

All games followed the iterated reference game framework, where a set of target images are described repeatedly over a series of blocks. Our games used 12 tangram images, previously used in Clark & Wilkes-Gibbs (1986) and Hawkins et al. (2020). Our games were run online, with crowd-sourced participants who were strangers to one another. The speaker knew which image was the target, and their goal was to describe it to the listeners over a chat interface so each listener could select the target. After all listeners had selected, players received feedback on the selections. The process repeated with the same speaker describing each of the 12 images to form one block. The games consisted of 6 blocks, for a total

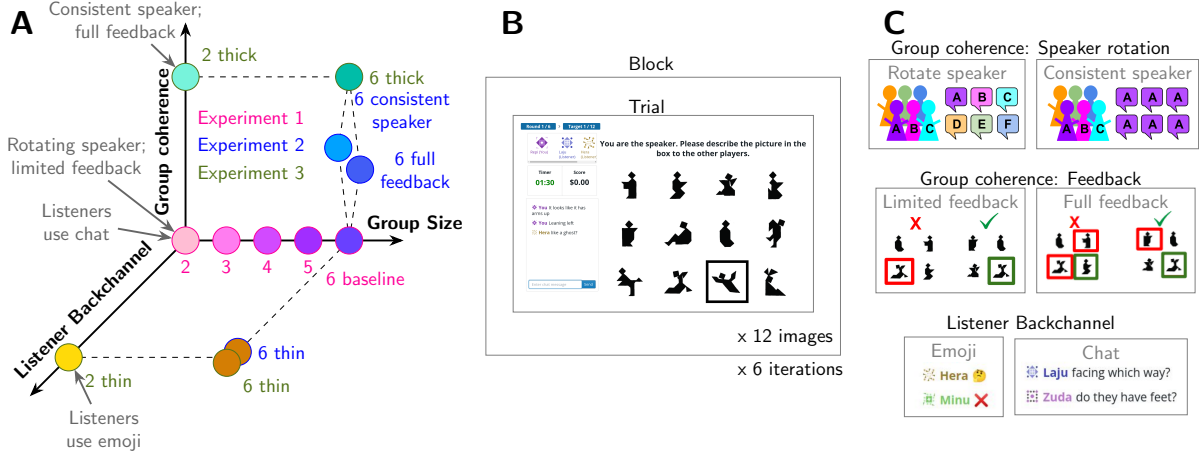


Figure 1: A: Diagram of the experimental space explored in these experiments. Experiments varied along 3 dimensions: Group size, group coherence, and listener backchannel. Each condition is shown as a dot. Experiment 1 (pink labels) varied group size from 2-6 players while holding group coherence and backchannel constant. Experiment 2 (blue labels) keep group size constant at 6 and varied the other dimensions. Relative to experiment 1, 6 consistent speaker and 6 full feedback each added one component of group coherence, and 6 thin reduced the backchannel. Experiment 3 (green labels) tested 4 corners of the space, crossing group size (2 or 6 players) with thin games (low coherence, low backchannel) or thick games (high coherence, high backchannel). B: Each trial a speaker described a target image to the listeners, and this process repeated for all 12 images to comprise a block, and the block repeated for a total of 6 iterations. C: Differences between conditions. See text for explanation.

of 72 trials, where each image was described 6 times over the course of the game (Figure 1B).

In experiment 1, we varied the size of the games continuously from 2 to 6 players while keeping other factors constant to explore what how performance varied across group size. For these conditions, the *speaker role rotated* after each block, so that all players had at least one turn as speaker. We gave *limited feedback* to listeners, showing each listener only whether their individual selection was correct or not, but not revealing others' selections or the right answer. Listeners had access to a *chat box* as their backchannel, so they could freely type questions and offer their own descriptions to the group.

In experiment 2, we explored varying different factors within 6 player games. We tried two different variations that we thought might improve group coherence and lead to better performance: having a *consistent speaker* rather than a rotating speaker, and separately, showing all the listeners *full feedback* on what each person in the group had selected and what the right answer had been. Additionally, we tested the role of listener contributions in establishing mutual understanding with a condition where listeners' backchannel was limited to four *emojis*. Listeners could send 4 discrete messages (green check, thinking face, red x, and laughing-crying face) to the chat. This limited backchannel allowed listeners to convey valence and level of comprehension, but not to contribute any referential content.

In experiment 3, we crossed the extremes of group size from experiment 1 (2 or 6 people) with extremes of group interactions from experiment 2. In the *thick* condition, we combined a consistent speaker with full feedback to create *high group coherence* and let listeners use the chat freely. In the *thin* condition, we repeated the emoji backchannel condition from experiment 2, which had *low group coherence* and an emoji backchannel. The 2-player thick game was similar to the condition used in Hawkins et al. (2020).

Group Performance

In dyadic reference games, listener accuracy remains high while the amount of referential language decreases dramatically over repetitions Hawkins et al. (2020). Listener accuracy measures how successful speakers are at communicating the target referent to the listeners. The combination of accuracy and reduction of speaker descriptions indicates that speaker and listeners have formed a shared conceptualization of the target that can be distilled into a shorter form, while still retaining the same level of informativity to listeners.

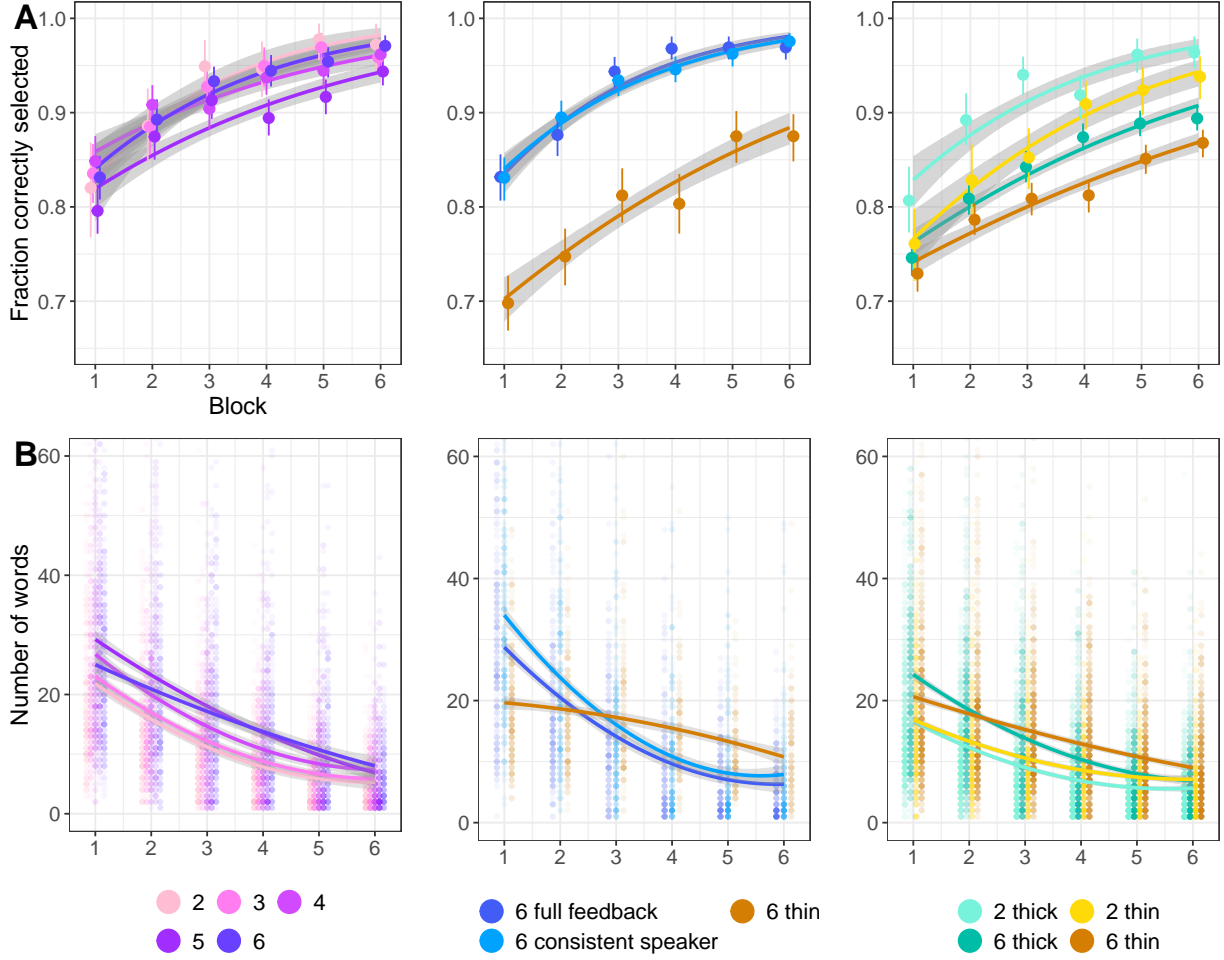


Figure 2: Behavioral results across all three experiments. A. Listener accuracy at selecting the target image. Dots are per condition, per block estimates with 95% bootstrapped CIs. Smooths are binomial fit lines. B. Number of words said by the speaker each trial. Faint dots represent individual trials from individual games. Smooths are quadratic fit lines. Y-axis is truncated, and a few outliers points are not visible.

Accuracy is high and increasing

In our experiments, listener accuracy rose over repetitions in all conditions and approached ceiling in most conditions; however, 6 player thin games were the least accurate (Figure 2A). We constructed logistic models to predict accuracy as a function of the manipulated variables, block, and their interactions.

Group size did not have a strong effect on initial accuracy or improvement rate ($\beta = -0.02$, 95% CrI = $[-0.05, 0.01]$) in the baseline condition (expt 1, $\beta = -0.07$, 95% CrI = $[-0.2, 0.05]$). 6 player thick games had lower initial accuracy (expt 3, $\beta = -0.64$, 95% CrI = $[-1.05, -0.25]$) and were slower to improve were slower to improve than 2-person thick games (expt 3, $\beta = -0.34$, 95% CrI = $[-0.43, -0.25]$).

Initial accuracy was somewhat higher for consistent speaker games ($\beta = 1.78$, 95% CrI = $[1.4, 2.19]$) and full feedback games ($\beta = 1.35$, 95% CrI = $[0.59, 2.06]$) than for 6 thin games (expt 2, $\beta = 0.88$, 95% CrI = $[0.64, 1.12]$), and improvement rates were higher for consistent speaker games ($\beta = 0.45$, 95% CrI = $[0.39, 0.52]$) and full feedback games ($\beta = 0.47$, 95% CrI = $[0.39, 0.54]$) than for 6 thin games (expt 2, $\beta = 0.23$, 95% CrI = $[0.19, 0.28]$). Thick and thin 2-player games were not reliably different in their initial accuracy (expt 3, $\beta = -0.36$, 95% CrI = $[-0.78, 0.05]$) or improvement rate (expt 3, $\beta = -0.07$, 95% CrI = $[-0.18, 0.04]$). There was no interaction between thinness and group size on initial condition ($\beta = 0.31$, 95% CrI = $[-0.22, 0.87]$) or improvement rate ($\beta = 0.07$, 95% CrI = $[-0.05, 0.19]$). Across experiments, the high and increasing levels of accuracy indicate that across all of these conditions, participants are able to succeed in communicating about the images.

Speaker utterances shorten over repetition

Speakers in larger games were more verbose than speakers in smaller games, and in some cases, these speakers showed sharper reduction from initial wordiness to eventual concision (Figure 2B). We constructed linear models of the number of words a speaker said each trial as a function of condition and block.

In the initial round, speakers in larger groups said more than speakers in smaller groups (expt 1, $\beta = 1.66$, 95% CrI = [0.66, 2.61]; expt 3, $\beta = 7.41$, 95% CrI = [3.57, 11.18]). The rate of reduction was relatively similar across group size for baseline games (expt 1, $\beta = -0.1$, 95% CrI = [-0.36, 0.17]), but 6-player games reduced faster than 2-player games in thick games ($\beta = -1.21$, 95% CrI = [-2.06, -0.3]). Taken together, smaller games on average said less at each stage of the game than larger games, although the gap closed in later rounds. Thin 2-player games were similar to thick 2-player games in initial verbosity (expt 3, $\beta = 0.63$, 95% CrI = [-3.18, 4.73]) and in reduction rate (expt 3, $\beta = 0.32$, 95% CrI = [-0.65, 1.24]). There was no effect of the interaction between thinness and game size on initial verbosity ($\beta = -2.11$, 95% CrI = [-7.41, 2.98]) or reduction rate ($\beta = 0.64$, 95% CrI = [-0.61, 1.89]).

The main effect of being one block later was -3.37, (95%CrI = [-4.54, -2.24]) words per trial in Experiment 1; -5.39, (95%CrI = [-6.46, -4.31]) words per trial for 6 consistent speaker; -4.68, (95%CrI = [-5.88, -3.52]) words for 6 full feedback, -2.15, (95%CrI = [-3.44, -1.12]) words for 6 thin; and $\beta = -2.29$, 95% CrI = [-2.95, -1.6] for 2-player thick. Overall, all conditions exhibited reduction, with some variation in the slope of reduction.

Listener utterances are more common in larger games

Listeners talked much less than speakers, and listener contributions were concentrated in the earlier trials. The more listeners, the more likely it was that some listener talked, and the more listeners said if they talked (TODO CROSSREF SUPPLEMENT IMAGE). We fit a logistic model predicting if listeners said anything, and for trials where they did, we fit a linear model predicting how many total words listeners said. The number of trials where any listener said anything related to the image was higher in larger groups (expt 1, $\beta = 0.78$, 95% CrI = [0.58, 0.98]) and declined across blocks (expt 1, $\beta = -0.8$, 95% CrI = [-0.97, -0.63]). When referential language was produced, larger groups produced more language (expt 1, $\beta = 2.12$, 95% CrI = [1.03, 3.12]), but the difference in group size closed in later blocks (expt 1, $\beta = -0.41$, 95% CrI = [-0.72, -0.09]). This pattern is consistent with early listener involvement in establishing a common conceptualization by asking questions and offering alternative descriptions. Once a shared idea is in place, listener descriptions are rarer and more perfunctory. Emoji use is not directly comparable to referential language, but similar trends occurred in the thin games: emoji use was more common in the 6-player games than the 2-player games and decreased over the course of the game (SUPP FIGURE).

Interim summary that TODO needs a better name

According to these performance and efficiency metrics, larger games are similar to smaller games, except with more talking, especially early in the games. Larger groups seem to generally take the time to elaborate descriptions in order for most listeners to understand, especially when listeners can ask specific clarifying questions. This leads to more talking in early rounds for larger games, sometimes followed by sharp reduction once a shared conceptualization is agreed upon.

Linguistic Content

Partner-specific reduction is characterized by high accuracy and shortening utterances, but the key phenomena of interest is that pairs are forming joint conventions about how to refer to particular images. How well groups are able to do this may vary based on game condition: we might predict that it's easier for two people to agree on a label, but harder when the labels are coming from 6 different individuals who may have different conceptualizations. Convergence to a shared name might be faster if only one person provides the labels, since they can match their own labels. Feedback and listener contributions may help a group get on the same page about what descriptions go with what images, letting a speaker reduce to shorthand expressions. While initially many groups may overlap with descriptions that include descriptions of the shapes or body parts in the image, their descriptions are predicted to become increasingly dissimilar as these descriptive portions drop out, leaving just the conventionalized nicknames that are group-dependent.

To assess the linguistic patterns of speaker descriptions, we examine the *semantic similarity* of descriptions

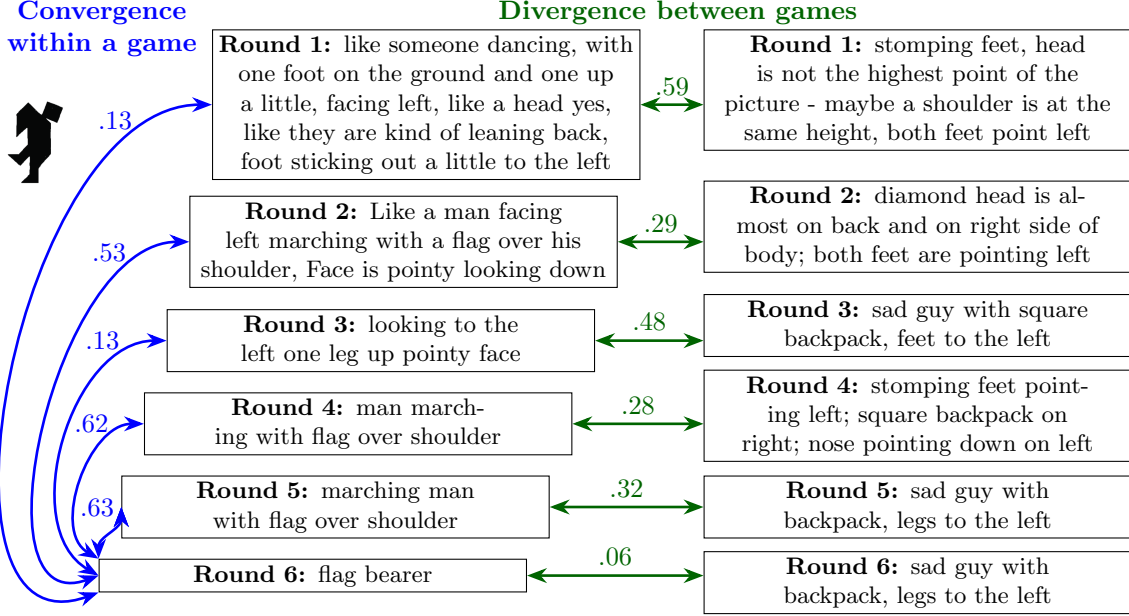


Figure 3: Example utterances describing the shown tangram figure produced by two 3-player games in Experiment 1. To measure convergence within a game (blue), we measured the cosine similarity between SBERT embeddings of descriptions and the embedding of the round 6 utterance (taken to be the convention). Higher cosine similarity indicates more similar meaning. To measure divergence between games (green), we measured the similarity between utterances from the same round across games.

within and across games. We quantified description similarity by concatenating speaker messages together within a trial and embedding this description into a high-dimensional vector space using SBERT (Reimers & Gurevych 2019). Then, we compared the similarity between pairs of utterances by taking the cosine similarity between their embeddings. Figure 3 shows an example of concatenated utterances and their cosine similarities. We looked at the similarity between utterances within a game compared to the final description of that image as a measure of convergence to a shared nickname. We looked at similarity between utterances across games for the same image in the same block as a measure of divergence between games, predicted to occur as each develops a different nickname.

Convergence within groups

Across conditions, speaker descriptions increased in semantic similarity to the final description over repetition; convergence was fastest in smaller and higher coherence groups, and was least strong in the 6-player thin condition (Figure 4A). We modeled semantic convergence by looking at the similarity between a block 1-5 utterance and the corresponding block 6 utterance as a function of the earlier block number and condition.

Smaller groups converged more strongly than larger groups. In the baseline condition, the similarity of the first utterance to the last utterance was invariant across group size (expt 1 $\beta = -0.008$, 95% CrI = $[-0.021, 0.005]$), but smaller groups converged faster (expt 1, $\beta = -0.008$, 95% CrI = $[-0.011, -0.005]$). In thick games, 6-player games had initially lower similarity (expt 3, $\beta = -0.069$, 95% CrI = $[-0.113, -0.025]$), but caught up to 2-player games (expt 3, $\beta = 0.009$, 95% CrI = $[0.001, 0.017]$). In both experiments, smaller games were always on average more similar to their latter utterances than larger games.

Thicker games converged more strongly than thinner groups. Convergence was especially rapid for the consistent speaker condition where all utterances came from the same person (expt 2, $\beta = 0.086$, 95% CrI = $[0.078, 0.094]$). Convergence was stronger for full-feedback (expt 2, $\beta = 0.062$, 95% CrI = $[0.051, 0.072]$) than 6 thin (expt 2, $\beta = 0.023$, 95% CrI = $[0.013, 0.033]$), and for 2-player thick games than 2-player thin games (expt 3, $\beta = -0.025$, 95% CrI = $[-0.033, -0.017]$).

Convergence was especially slow and weak for thin 6-player games (expt3, $\beta = -0.035$, 95% CrI = $[-0.047, -0.025]$). Convergence towards the last utterance was driven by cumulative increasing similarity between pairs of utterances in adjacent blocks (see TODO SUPPLEMENT). Overall, descriptions were

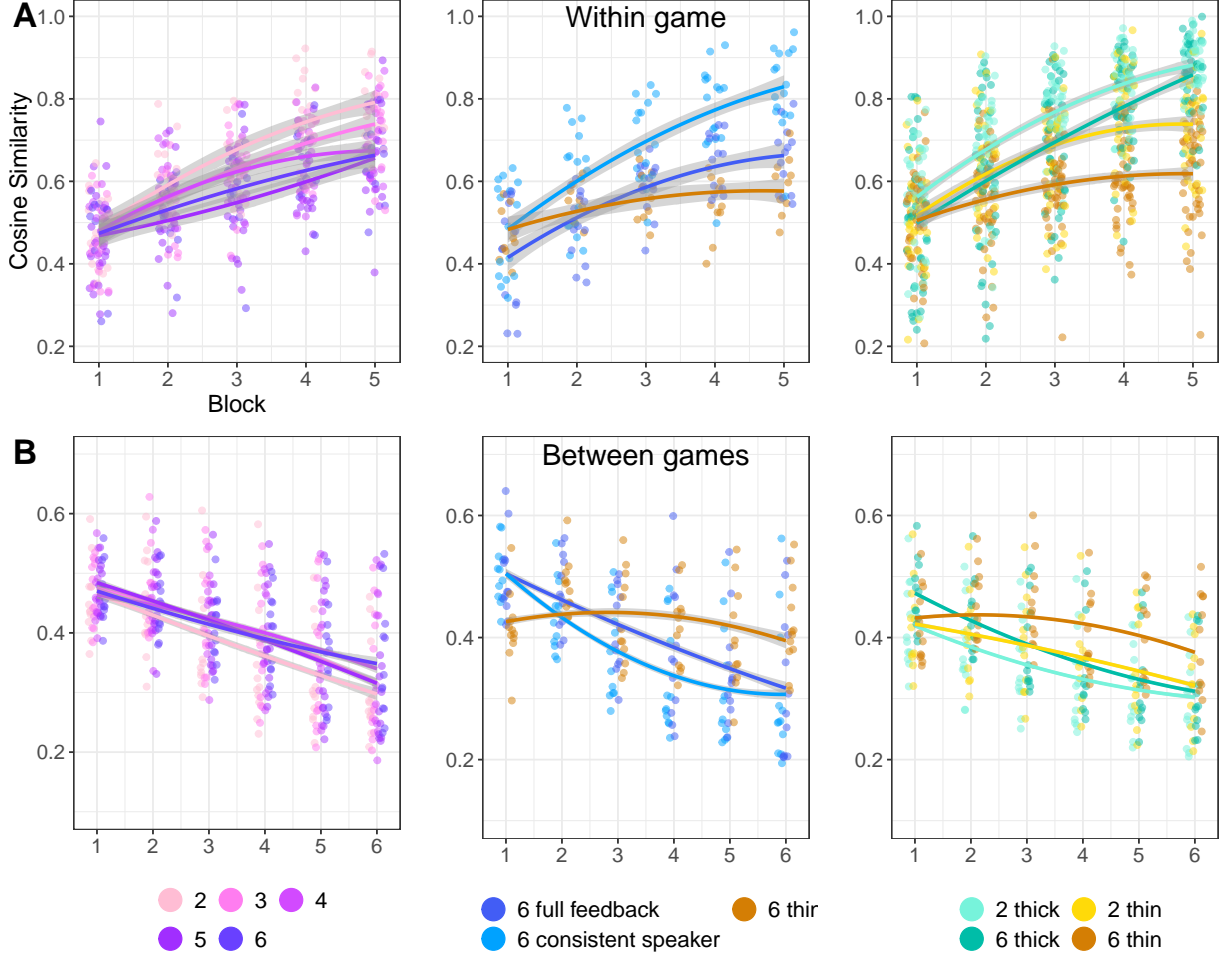


Figure 4: Language similarity results measured with pairwise cosine similarity between embeddings of two utterances. A. Convergence of utterances within games as measured by similarity between an utterance from block 1-5 to the block 6 utterance in the same game for the same image. Dots are per-game averages, smooths are quadratic. B. Divergence of utterances across games as measured by the similarity between an utterances and utterances produced for the same image by different groups in the same block. Dots are per-image averages, smooths are quadratic.

more similar if provided by the same person, if fewer people were in the game, and if listeners could contribute via a text channel.

Divergence between groups

Over repetitions, speaker descriptions diverged from descriptions used in other groups: divergence was fastest in groups with thick communication channels, while the 6-player thin condition games barely diverged at all (Figure 4B). We modeled semantic divergence by looking at the similarity between a pair of utterances for the same image from the same block across different games as a function of the block number and condition.

Game size had inconsistent effects on divergence slope: in baseline games, smaller games diverged faster than larger games (expt 1, $\beta = 0.001$, 95% CrI = [0.001, 0.002]) from the same starting point (expt 1, $\beta = 0.002$, 95% CrI = [0, 0.004]). In thick games, 6-player games diverged faster than 2-player games (expt 3, $\beta = -0.008$, 95% CrI = [-0.01, -0.007]) but from initially higher cross-game similarity (expt 3, $\beta = 0.051$, 95% CrI = [0.047, 0.055]). In both cases, smaller games were less similar to each other on average.

Thicker games tended to diverge faster than thinner games. Both consistent speaker (expt 2, $\beta = -0.041$, 95% CrI = [-0.043, -0.039]) and full feedback games (expt 2, $\beta = -0.038$, 95% CrI = [-0.04, -0.035]) diverged faster than thin games. Thick 2-person games diverged faster than thin

2-person games (expt 3, $\beta = 0.004$, 95% CrI = [0.002, 0.005]). Most noticeably, the 6 player thin games diverged much more slowly than the other conditions (expt 3, $\beta = 0.017$, 95% CrI = [0.015, 0.019]).

Overall, thicker communication channels seem to be necessary for fast divergence, although 2-player games can diverge, albeit more slowly, even with a thin channel. Similarly to the convergence data, the 6 player thin condition is an outlier that barely follows trends present across other conditions.

Different patterns of performance and content

Theoretical approaches treat reduction in referring expression length as a consequence of partner-specific convention formation (Clark & Wilkes-Gibbs 1986, Brennan & Clark 1996, Yoon & Brown-Schmidt 2014, 2018), but in our current work, reduction and semantic measures of partner-specificity pattern differently in the 6 player thin condition. The 6-player thin games show much less divergence between games and convergence within games, even compared to the 2 thin and 6 thick conditions, but 6 player thin games showed smaller (and statistically inconclusive) differences to 6 thick games for accuracy and reduction. This gap between the group performance measures and the semantic measures raises the possibility that it is possible to become more concise (and more accurate) without developing group-specific nicknames, but instead perhaps relying on group priors and reducing the amount of detail (Guilbeault et al. 2021). This gap highlights the need to use measures of the type of language (and not just amount of language) when looking for convention-formation phenomena.

General Discussion

Communication often occurs in multi-party settings, but research on referential communication often does not. In dyadic work, iterated reference games have been used to establish a phenomena of reduction over repeated reference, characterized by speaker-listener pairs creating short nicknames that they mutually understand, but which are not shared by other groups. In this work, we asked how this process of reference formation unfolds under varying interaction structures.

Across 3 online experiments and 11 experimental conditions, we varied game features including group size, form of listener backchannel, and degree of group coherence. All conditions showed the hallmarks of reduction: increasing accuracy, reduction in speaker utterances, semantic convergence within games, and differentiation of descriptions between groups. Even with larger groups and more constrained means of communication, reduction still occurs.

However, while results are directionally the same across conditions, the interaction structure of a group substantially affects how rapidly groups develop partner-specific conventions. Smaller groups and games with thicker communication channels converged faster and more robustly than games that were larger or had thinner communication channels. These factors add together to form the overall group experience. The differences between the 6 player thin condition and both the 2 player thin condition and other 6 player conditions point to an interaction: 2-player games can cope with limited feedback mechanisms, but 6-player games suffer without access to more feedback. Group dynamics differ depending on group size, and larger groups may be more sensitive to other factors affecting interaction structure. Multi-player groups thus make for a richer and more sensitive environment to study communication phenomena applicable to both pairs and small groups.

Just within the general framework of iterated reference, there is a high dimensional feature space of possible experiments. We sampled only a few points along a few dimensions in the space that felt salient. In our experiment 3, we grouped some factors together in order to have more games in each condition: a fully factorial design would have been too expensive to power adequately. Future work could sample other points in the experimental space, perhaps exploring the effects of different target images, or groups of people with real-life prior connections.

We cannot make claims about causal mechanisms between how experimental set-ups such as group size resulted in different outcomes: for instance, there are many differences between being in a 2-person group versus a 6-person group that could lead to the different outcomes. In a dyad, speakers can tailor their utterances to the one listener, but in large groups, speakers must balance the competing needs of different listeners (Schober & Clark 1989, Tolins & Fox Tree 2016). These effects likely vary by both the knowledge state of and communication channels available to the listeners (Horton & Gerrig 2002, Horton & Gerrig 2005, Fox Tree & Clark 2013). Further work digging into the language used and the interactions

between participants might unearth plausible mechanisms for how differences in group size and interaction structure influence outcomes, and this in turn could then point towards future experimental conditions.

Methods

For all experiments, we used Empirica (Almaatouq et al. 2020) to create real-time multi-player iterated reference games. In each game, one of the players started as the speaker who saw an array of tangrams with one highlighted and communicated which figure to click to the other players (listeners). After the speaker had identified each of the 12 images in turn, the process repeated with the same images, but a total of 6 blocks (72 trials). We recorded what participants said in the chat, as well as who selected what image and how long they took to make their selections.

These experiments were designed sequentially and pre-registered individually.¹ We followed the analysis plan, although additional analyses were added to early experiments that were only pre-registered in later experiments. Results from some pre-registered models were omitted from the main text, but are shown in the supplement TODO.

Participants

Participants were recruited using the Prolific platform, and all participants self-reported as fluent native English speakers on Prolific’s demographic prescreen. Participants each took part in only one experiment. Experiment 1 took place between May and July 2021, experiment 2 between March and August 2022, and experiment 3 in October 2022. As games varied in length depending on the number of participants, we paid participants based on group size, with the goal of a \$10 hourly rate. 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. When one player had the speaker role for the entirety of a 6-player game, they gained an addition \$2 bonus. Across all games, each participant could earn up to \$2.88 in performance bonuses. A total of 1319 people participated across the 3 experiments, for roughly 20 games in each condition in experiments 1 and 2 and 40 games per condition in experiment 3. A breakdown of number of games and participants in each condition is shown in TODO supplement.

Materials

We used the 12 tangram images used by Hawkins et al. (2020) and Clark & Wilkes-Gibbs (1986). These images were displayed in a grid with order randomized for each participant (thus descriptions such as “top left” were ineffective as the image might be in a different place on the speaker’s and listeners’ screens). The same images were used every block.

Procedure

The experimental procedure was very similar across the three experiments. We first describe the procedure used in experiment 1 and then describe the differences in later experiments.

Experiment 1

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 partner(s) were ready.

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. There was no signal to the speaker or other listeners about who had already made a selection.

Once all listeners had selected (or a 3-minute timer ran out), participants were given feedback. Listeners 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, but listeners did not.

¹Experiment 1: <https://osf.io/cn9f4> for the 2-4 player groups, and <https://osf.io/rpz67> for the 5-6 player data run later. Experiment 2: consistent speaker at <https://osf.io/f9xyd>, full feedback at <https://osf.io/j5zbn>, and thin at <https://osf.io/k5f4t>. Experiment 3: <https://osf.io/untzy>

Listeners got 4 points for each correct answer; the speaker got points equal to the average of the listeners' points. These points translated into 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. Rotating the speaker was chosen to keep participants more equally engaged (the speaker role is more work), and to give a more robust test for reduction and convention.

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

Differences in experiment 2

Experiment 2 consisted of three different variations on Experiment 1, all conducted in 6 player games. Each of these conditions differed from the experiment 1 baseline in one way. The consistent speaker condition differed only in that one person was designated the speaker for the entire game, rather than having the speaker role rotate. The full feedback condition differed from experiment 1 in that all participants were shown what each person had selected and what the right answer was; listeners still saw text saying whether they individually were right or wrong. This was similar to some dyadic work, such as Hawkins et al. (2020) where listeners were shown what the right answer was during feedback. For the thin condition, we altered the chatbox interface for listeners. Instead of a textbox, listeners had 4 buttons, each of which sent a different emoji to the chat. Listeners were given suggested meanings for the 4 emojis during instructions. They could send the emojis as often as desired, for instance, initially indicating confusion, and later indicating understanding. In addition, we added notifications that appeared in the chat box saying when a player had made a selection.

Differences in experiment 3

The thin channel condition in experiment 3 was the same as the thin condition in experiment 2, above. The thick condition combined the two group coherency enhancing variations from experiment 2: one person was the designated speaker throughout, and the feedback participants received included the right answer and what each player had selected. Across both conditions in experiment 3, notifications were sent to the chat to indicate when a participant had made a selection.

Data pre-processing and exclusions

Participants could use the chat box freely, which meant that the chat transcript contained some non-referential language. The first author skimmed the chat transcripts, tagging utterances that did not refer to the current tangram. These were primarily pleasantries ("Hello"), meta-commentary about how well the task was going, and confirmations or denials ("ok", "got it", "yes", "no"). We excluded 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.

In experiments 1 and 2, games did not start if there were not enough participants and ended if any participant disconnected. In experiment 3, games started after a waiting period even if they were not full and continued even after a participant disconnected (with speaker role reassigned if necessary), unless the game would drop below 2 players. The distribution of players in these 6* player games is at TODO SUPPLEMENT! The realities of online recruitment and disconnection meant that the number of games varied between conditions. We excluded incomplete blocks from analyses, but included complete blocks from partial games (See SUPPLEMENT TABLE for counts).

When skimming transcripts to tag non-referential utterances, we noticed that one game in the 6-player thick game had a speaker who did not give any sort of coherent descriptions, even with substantial listener prompting. We excluded this game from analyses.

Modelling strategy

In experiment 3, some of the 6 player games did not have 6 players for the entire game. We do not model this, as it is unclear at what point in the game group size is most relevant. We note that this is a conservative choice that will underestimate differences between 2 player and (genuine) 6 player games, by labelling some smaller groups as 6 player.

We ran all models in BRMS (Bürkner 2018) with weakly regularizing priors. We were often unable to fit the full mixed effects structure that we had pre-registered in a reasonable amount of time, so we included what hierarchical effects were reasonable. (All model results and priors and formulae are reported in TODO supplement). Accuracy models were run as logistic models with normal(0,1) priors for both betas and sd. Reduction models were run as linear models with an intercept prior of normal(12,20), a beta prior of normal(0,10), an sd prior of normal(0,5) and a correlation prior of lkj(1). For all of the models of sbert similarity, we used linear models with the priors normal(.5,.2) for intercept, normal(0,.1) for beta, and normal(0,.05) for sd.

References

- Almaatouq A, Becker J, Houghton JP, Paton N, Watts DJ, Whiting ME (2020) [Empirica: A virtual lab for high-throughput macro-level experiments](#). *ArXiv200611398 Cs*
- Almaatouq A, Griffiths TL, Suchow JW, Whiting ME, Evans J, Watts DJ (2022) Beyond Playing 20 Questions with Nature: Integrative Experiment Design in the Social and Behavioral Sciences. *Behavioral and Brain Sciences*:1–55. doi:[10.1017/S0140525X22002874](#)
- Branigan H (2006) Perspectives on multi-party dialogue. *Research on Language and Computation* **4**:153–177
- Brennan SE, Clark HH (1996) Conceptual Pacts and Lexical Choice in Conversation. :12
- Bürkner P-C (2018) Advanced bayesian multilevel modeling with the r package brms. *The R Journal* **10**:395–411
- Carletta J, Garrod S, Fraser-Krauss H (1998) Placement of Authority and Communication Patterns in Workplace Groups: The Consequences for Innovation. *Small Group Research* **29**:531–559. doi:[10.1177/1046496498295001](#)
- Clark HH, Wilkes-Gibbs D (1986) [Referring as a collaborative process](#). *Cognition*
- Cohn-Gordon Reuben, Levy R, Bergen L (2019) The pragmatics of multiparty communication.
- Fay N, Garrod S, Carletta J (2000) Group Discussion as Interactive Dialogue or as Serial Monologue: The Influence of Group Size. *Psychol Sci* **11**:481–486. doi:[10.1111/1467-9280.00292](#)
- Fox Tree JE, Clark NB (2013) Communicative Effectiveness of Written Versus Spoken Feedback. *Discourse Processes* **50**:339–359. doi:[10.1080/0163853X.2013.797241](#)
- Ginzburg J, Fernandez R (2005) Action at a distance: The difference between dialogue and multilogue. *Proceedings of DIALOR*:9
- Guilbeault D, Baronchelli A, Centola D (2021) Experimental evidence for scale-induced category convergence across populations. *Nat Commun* **12**:327. doi:[10.1038/s41467-020-20037-y](#)
- Haber J, Baumgärtner T, Takmaz E, Gelderloos L, Bruni E, Fernández R (2019) The PhotoBook Dataset: Building Common Ground through Visually-Grounded Dialogue. In: *Proc. 57th Annu. Meet. Assoc. Comput. Linguist.* Association for Computational Linguistics, Florence, Italy, p 1895–1910. Available from: <https://www.aclweb.org/anthology/P19-1184> [Last accessed 1 February 2022]. doi:[10.18653/v1/P19-1184](#)
- Hawkins RD, Frank MC, Goodman ND (2020) [Characterizing the dynamics of learning in repeated reference games](#). *ArXiv191207199 Cs*
- Horton WS, Gerrig RJ (2002) Speakers’ experiences and audience design: Knowing when and knowing how to adjust utterances to addressees. *Journal of Memory and Language*:18
- Horton WS, Gerrig RJ (2005) The impact of memory demands on audience design during language production. *Cognition* **96**:127–142. doi:[10.1016/j.cognition.2004.07.001](#)
- Krauss RM, Weinheimer S (1964) Changes in reference phrases as a function of frequency of usage in social interaction: A preliminary study. *Psychon Sci* **1**:113–114. doi:[10.3758/BF03342817](#)
- Krauss RM, Weinheimer S (1966) Concurrent feedback, confirmation, and the encoding of referents in verbal communication. *Journal of Personality and Social Psychology* **4**:343–346. doi:[10.1037/h0023705](#)
- Metzing C, Brennan SE (2003) When conceptual pacts are broken: Partner-specific effects on the comprehension of referring expressions. *Journal of Memory and Language* **49**:201–213. doi:[10.1016/S0749-596X\(03\)00028-7](#)

- Murthy SK, Griffiths TL, Hawkins RD (2022) Shades of confusion: Lexical uncertainty modulates ad hoc coordination in an interactive communication task. doi:[10.48550/arXiv.2105.06546](https://doi.org/10.48550/arXiv.2105.06546)
- Reimers N, Gurevych I (2019) Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. doi:[10.48550/arXiv.1908.10084](https://doi.org/10.48550/arXiv.1908.10084)
- Rogers SL, Fay N, Maybery M (2013) Audience Design through Social Interaction during Group Discussion. *PLOS ONE* **8**:e57211. doi:[10.1371/journal.pone.0057211](https://doi.org/10.1371/journal.pone.0057211)
- Schober MF, Clark HH (1989) Understanding by addressees and overhearers. *Cognitive Psychology* **21**:211–232. doi:[10.1016/0010-0285\(89\)90008-X](https://doi.org/10.1016/0010-0285(89)90008-X)
- Tolins J, Fox Tree JE (2016) Overhearers Use Addressee Backchannels in Dialog Comprehension. *Cogn Sci* **40**:1412–1434. doi:[10.1111/cogs.12278](https://doi.org/10.1111/cogs.12278)
- Traum D (2004) Issues in Multiparty Dialogues. In: Dignum F (ed) *Advances in Agent Communication*. Springer Berlin Heidelberg, Berlin, Heidelberg, p 201–211. Available from: http://link.springer.com/10.1007/978-3-540-24608-4_12 [Last accessed 1 February 2022]. doi:[10.1007/978-3-540-24608-4_12](https://doi.org/10.1007/978-3-540-24608-4_12)
- Weber RA, Camerer CF (2003) Cultural Conflict and Merger Failure: An Experimental Approach. *Manag Sci* **49**:16
- Yoon SO, Brown-Schmidt S (2014) Adjusting conceptual pacts in three-party conversation. *Journal of Experimental Psychology: Learning, Memory, and Cognition* **40**:919–937. doi:[10.1037/a0036161](https://doi.org/10.1037/a0036161)
- Yoon SO, Brown-Schmidt S (2018) Aim Low: Mechanisms of Audience Design in Multiparty Conversation. *Discourse Processes* **55**:566–592. doi:[10.1080/0163853X.2017.1286225](https://doi.org/10.1080/0163853X.2017.1286225)
- Yoon SO, Brown-Schmidt S (2019) Audience Design in Multiparty Conversation. *Cogn Sci* **43**:e12774. doi:[10.1111/cogs.12774](https://doi.org/10.1111/cogs.12774)