

# TODO title

Veronica Boyce<sup>1,\*</sup>, Robert Hawkins<sup>1</sup>, Noah D. Goodman<sup>1</sup>, Michael C. Frank<sup>1</sup>

<sup>1</sup>Stanford University

## Abstract

*This is an abstract in italics.*

This is the second paragraph not in italics.

## Keywords

One keyword; Yet another keyword

## 1 Introduction

- communication is important and takes place in non-dyadic situations
- one interesting thing about comm is ad-hoc expressions, adaptation to speaker/ listener pairs
- prior work has focused a lot on dyads: what happens with more people? Do theories predict anything?
- groups also vary in compositional type aspects so mimics of these might matter
- do we want to set up other theoretical frames? do I need to read more about theory??

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. Communicative environments range in size from one-on-one dialogue to broadcast communication to large groups, but the goal of efficient communication is shared across these (Traum 2004, branigan2006?, ginzburg2005?). Shared referring expressions are a necessity for efficient communication; a thing or an idea needs some sort of name that the interlocutors will jointly understand. In many cases, there are widely shared conventionalized expressions for objects or ideas, but in other cases, spontaneous ad-hoc expressions must be invented.

The formation of these new reference expressions is well-studied in dyadic contexts and has been a case study for efficient communication more broadly. But these 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 larger groups.

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, krauss-ConcurrentFeedbackConfirmation1966?). Two participants see the same set of tangram 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 with the same images over a series of blocks. Early descriptions are long and make reference to multiple features in the figure, but in later iterations, shorthand conventional names for each figure emerge; this shortening of utterances is called ‘reduction’.

---

\*Corresponding author. Email: [vboyce@stanford.edu](mailto:vboyce@stanford.edu)

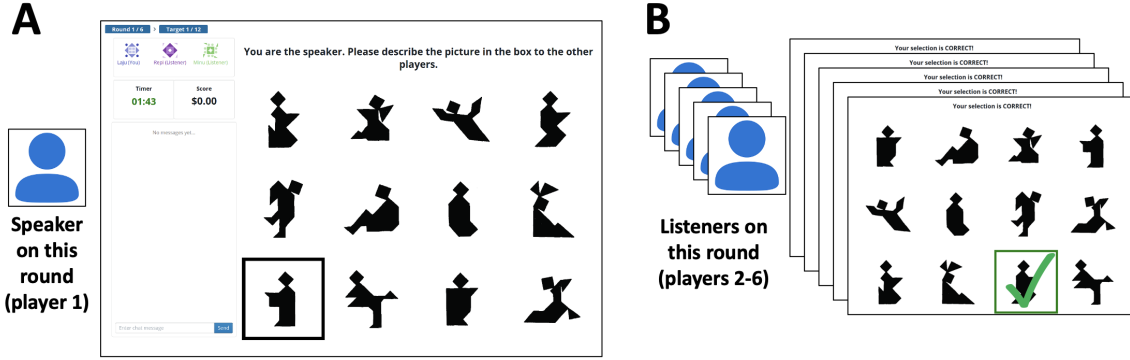


Figure 1: All participants saw all 12 tangram images. (A) Speaker’s view during selection phase. (B) During the feedback stage, speakers saw what figure each person chose, but listeners only learned if their selection was correct or incorrect. Listeners were not shown what other listeners chose.

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 Hawkins et al. (2020), 83 pairs completed a similar iterated reference experiment where they communicated via a chat box. Speakers reduced their utterances, producing fewer words per image in later blocks than in earlier blocks, in line with results from face-to-face, oral paradigms.<sup>1</sup>

How does this process proceed in multi-party communication? 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 (fox-tree2013?, horton2002?, horton2005?). Prior work has focused on manipulating knowledge states by adding new listeners to established groups.

In this context, one approach for speakers is to ‘aim low’ and produce utterances tailored to the least knowledgeable listener (yoonAimLowMechanisms2018a?). For instance, in Yoon & Brown-Schmidt (2014), speakers developed conventions with one listener but then used longer descriptions with a new listener. Another strategy for speakers is to integrate across listeners and balance efficiency with informativeness by ‘aiming in the middle’. In Yoon & Brown-Schmidt (2019), speakers communicating to a mixed group of 3 experienced listeners and 1 naive listener used shorter utterances and made fewer accommodations than they did in groups with a greater fraction of naive listeners. Both of these strategies predict that larger groups will be slower to converge than smaller groups.

Disagreements about how to conceptualize referents can also slow groups down. In (weberCulturalConflictMerger2003?), pairs of participants played a reference game with the same image sets before a listener switched groups and joined a different pair, making a group of three. The addition of the new listener slowed both listeners down for multiple rounds. When a listener switched groups, they brought preconceptions about how the pictures should be described which conflicted with how the speaker was used to describing the images. This result predicts that, with more perspectives in play, larger groups may have more difficulty agreeing on common conceptualizations.

In general, listeners expect speakers to maintain conventions and stick to descriptions that were similar to successful descriptions. However, listeners were not surprised to hear different descriptions of a familiar object if it came from a new speaker who had just entered the room (Metzing & Brennan 2003). It’s unclear what this finding predicts about new speakers who are present as fellow listeners during prior blocks – will listeners expect them to maintain conventions?

<sup>1</sup>We use “speaker” and “listener” to refer to the roles describing and selecting targets, regardless of communication modality.

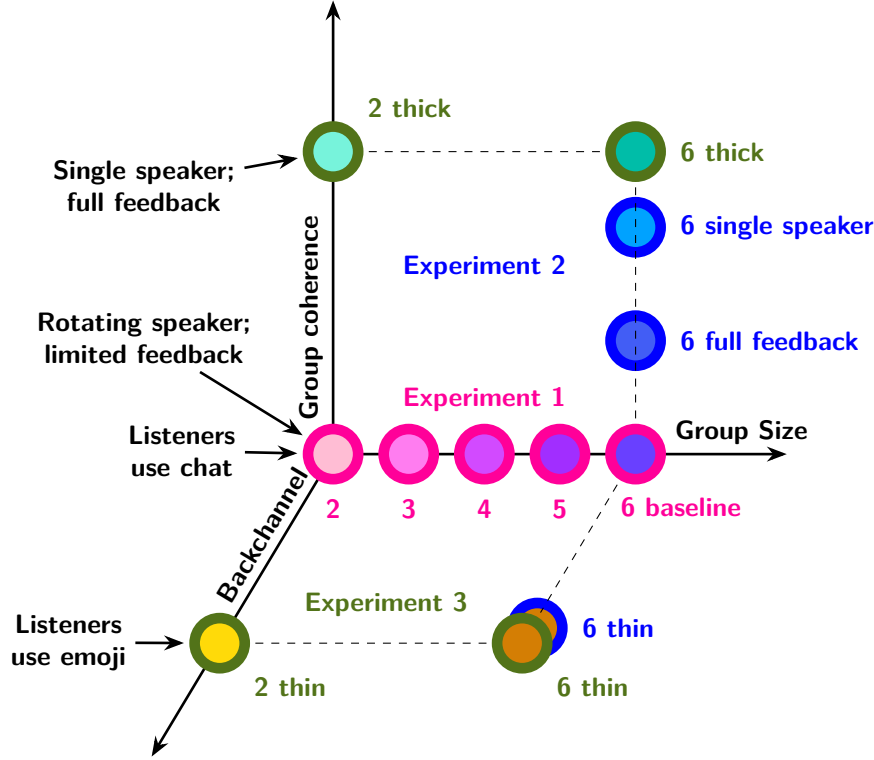


Figure 2: FOOBAR

## 2 Methods

### 2.1 fluff

For a first experiment, 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. This paradigm allowed us to test how the findings from dyadic iterated reference games extend to larger groups and were influenced by group size. Additionally, to more closely examine how conventionalized names develop within groups and differ between groups, we used sentence embeddings to quantify the similarities of descriptions within and between groups over time. Building on the methods of Hawkins et al. (2020), we used Empirica (Almaatouq et al. 2020) to create real-time multi-player reference games. In each game, one of the players started as the speaker who saw an array of tangrams with one highlighted (Figure 1A) and communicated which figure to click to the other players (listeners). After the speaker had identified each of the 12 images in turn, the speaker role rotated to another player and the process repeated with the same images. In total, there were 6 blocks, giving each player at least one chance to be the speaker. We recorded what participants said in the chat, as well as who selected what image and how long they took to make their selections.

After experiment 1, there were concerns about whether the results were dependent on any of the methods details. We were also interested in whether there were changes in knowledge or group communication that could support (or diminish) success, especially in the larger groups. To explore this, we ran 3 variants of the paradigm in experiment 1, varying different parameters. These were all run on small sets of 6 player groups.

TODO examples of prior literature that rotated or didn't rotate !!! Previous literature varied on whether one person had the speaker role the entire game, or whether the two participants alternated speaker and listener roles in some way. For experiment 1, we went with rotating the speaker every block: we thought this would be perceived as more fair by participants, and we thought it would be a more stringent test of convergence to convention. However, in smaller groups, the speaker role would rotate back around, meaning that the speaker in later blocks would be more experienced in small games than in large games where speakers were usually first time speakers. We cannot disentangle the effects of group size from the effects of speaker experience. To address whether continuity in speakers might contribute to better

performance, in experiment 2a we ran 6 player games with one person assigned to be the speaker the entire game.

Another point of variation in the literature is what sort of feedback is given. In particular, there is variation in whether participants were told what the right answer was when they were wrong (like in Hawkins et al. (2020)). In experiment 1, we opted for the minimum feedback, to make sure that participants had to negotiate and understand conventions, and could not just pick them up from the end of trial feedback. To see whether this mattered, in experiment 2b, we increased the amount of feedback, showing each participant what everyone had selected and what the right answer was.

Listener feedback to the speaker in the form of some backchannel, is implicated in the theoretical literature as being important to convention formation TODO CITES. We were interested in whether limiting this backchannel would be a way to push around the reduction pattern and significantly change it. We were worried that entirely removing the backchannel would make the game too unfun and encourage random clicking or quitting of the speaker was just talking into the void. As a compromise between limiting the backchannel and maintaining playability, we switched to giving the listeners a discrete backchannel consisting of 4 emojis, expressing whether they understood, whether they wanted more details, whether they were totally confused, or whether they were amused (TODO figure with emojis and their meanings). These options were based on the common valences expressed in the free chat of listeners from experiment 1. Thus, in 2c, listeners could not contribute questions or content to the discourse, but the speaker still had some feedback on whether the listeners were understanding.

Based on the results of experiment 2, it seemed that speaker continuity helped groups, and that more feedback maybe helped groups, while limiting the backchannel reduced accuracy and reduction. These results are tentative as we had a limited number of groups in each of these conditions. To more robustly demonstrate that the reduction phenomena were sensitive to these variations in group structure and coherency and to investigate how they interacted with group size, we conducted a better powered experiment. As these experiments are expensive to run, we did not do a full-factorial design; instead, we collapsed the three sources of variation into a “thick-channel” condition, combining single speaker, high feedback, and text backchannel, the settings that give richer feedback and knowledge, and a “thin-channel” condition, which was the opposite (same condition as 2c above). We crossed these two conditions with group size: either 2 or 6 player groups, picking the extreme values from experiment 1.

We aimed for 40 games in each of these 4 cells. We had previously experienced problems with attrition in 6-player games, as when one or more participants quit, the game discontinued on everyone, leading to fewer games with data from the entire game. To mitigate this, we reprogrammed the game to instead try to keep the game going with the participants who were left (and to start games even if they hadn’t filled completely).

## 2.2 Methods

Pre-regs: 1: <https://osf.io/cn9f4> for the 2-4 player groups, and <https://osf.io/rpz67> for the 5-6 player data run later 2a: <https://osf.io/f9xyd> 2b: <https://osf.io/j5zbm> 2c: <https://osf.io/k5f4t> 3: <https://osf.io/untzy>

## 2.3 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 early up to \$2.88 in performance bonuses. A total of 1319 people participated across the 3 experiments.

Table 1: Summary of differences in experiments. Game size refers to the number of players per game. Speaker refers to whether there was one speaker the whole game or whether the speaker role rotated every block. Feedback is whether listeners saw only whether they were right or wrong or whether the additionally saw what other listeners had selected and what the correct answer was. Listener chat refers to whether listeners could type freely in the chat or only communicate by pressing buttons to send four emojis to the chat. Continue games refers to whether games could continue (or start) with fewer than the requisite number of players; this was not intended to be a consequential manipulation, but was done to prevent games from ending if one player dropped out (an issue that was causing data loss in 6 player games).

Experiment	Game size	Speaker	Feedback	Listener chat	Continue games
1	2,3,4,5,6	rotating	self only	text	no
2a	6	one speaker	self only	text	no
2b	6	rotating	self, others, & correct	text	no
2c	6	rotating	self only	four emojis	no
3 thin	2,6	rotating	self only	four emojis	yes
3 thick	2,6	one speaker	self, others, & correct	text	yes

Table 2: The number of games in each experiment and condition. Complete games finished all 6 blocks; partial games ended early due to disconnections, but contributed at least one complete block of data. 6\* indicates that some games started with fewer than 6 players or continued with fewer than 6 players after participants disconnected.

Experiment	Players	Complete	Partial
1	2	15	4
1	3	18	2
1	4	19	2
1	5	17	3
1	6	12	6
2a	6	15	3
2b	6	13	4
2c	6	10	6
3: thin	2	35	3
3: thin	6*	44	0
3: thick	2	39	3
3: thick	6*	38	2

## 2.4 Materials

We used the 12 tangram images used by Hawkins et al. (2020) and Clark & Wilkes-Gibbs (1986) (see Figure 1). 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.

## 2.5 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.

### 2.5.1 Experiment 1

We implemented the experiment using Empirica, a Javascript-based platform for running real-time interactive experiments online (Almaatouq et al. 2020). From Prolific, participants were directed to our website where they navigated through a self-paced series of instruction pages explaining the game. Par-

ticipants 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.

Once the game started, participants saw screens like Figure 1A. 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 (Figure 1B). 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. This feedback regime is different from Hawkins et al. (2020) where listeners were shown what the right answer was during feedback. We made this change to prevent listeners from learning conventions purely as a memorized mapping between utterance and correct answer.

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.

### 2.5.2 Differences in experiment 2

Experiment 2 consisted of three different variations on Experiment 1, so we describe the differences from the Experiment 1 procedure. For the differences between games see TABLE TODO.

The only change for experiment 2a was that one person was designated the speaker for the entire game.

For experiment 2b, feedback was changed so that all participants in the feedback stage saw the speaker view showing who selected what and what the right answer was. Listeners still saw text saying whether they individually were correct or wrong.

For experiment 2c, 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.

### 2.5.3 Differences in experiment 3

The thin channel condition was the same as with experiment 2c, described above.

The thick channel condition was the same as experiment 1, except that one person was designated the speaker throughout, and all participants saw the speaker view for feedback.

TODO confirm. Across both conditions, messages were sent to the chat to indicate when a participant had made a selection.

The large change was behind the scenes, to handle games continuing after a participant had quit. Participants were identified as disconnected if their computer was not responding to the server for YY seconds (for instance if they closed the tab and did not reopen it quickly). If the person who disconnected was a listener, they were just removed, and they were skipped over by the speaker rotation (if applicable). If the person who disconnected was a speaker, that trial was discontinued as there was no way for listeners to get more information, and another person was assigned as speaker for the remainder of the block (or remainder of the game, depending on condition).

TODO figure out how this interacts with the full blocks thing!!! Maybe we should break the pre-reg? (in general?, only here?)

TODO the actual distribution of game sizes!!

TODO add pre-registrations everywhere

Note: 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.

## 2.6 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 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 or condition for experiments 1 and 2, and 40 games in each condition in experiment 3. In some cases, we ended up with fewer due to games not filling or participants disconnecting early (Table ??). We excluded incomplete blocks from analyses, but included complete blocks from partial games.

## 2.7 Modelling strategy

note that accuracy model deviates from pre-reg b/c I went overboard with the mixed effects model structure

All models were run in brms CITE with weakly regularizing priors DESCRIBE

# 3 Results

## 3.1 Behavioral results

The two key behavioral result outcomes were the accuracies of the listeners at selecting the target tangrams and the number of words produced by the speaker over the course of each trial.

Across all experiments, most individuals were accurate in their selections, with accuracy rising across blocks, but accuracy was noticeably lower, but still far above chance, in the 6-player games with an emoji backchannel (Figure ??). In experiment 1, participants were more accurate in later blocks (block: Est=0.44, CrI=[0.31, 0.58]), and there was no strong effect of group size on accuracy (numPlayers: Est=-0.07, CrI=[-0.2, 0.05]) or interaction between block and group size (block:numPlayers: Est=-0.02, CrI=[-0.05, 0.01]). In experiment 2, participants were again more accurate in later blocks (2a: block: Est=0.45, CrI=[0.39, 0.52], 2b: block: Est=0.47, CrI=[0.39, 0.54], 2c block: Est=0.23, CrI=[0.19, 0.28]). In experiment 3, participants are more accurate in later blocks (block: Est=0.41, CrI=[0.32, 0.5]). Participants are less accurate overall in the six player games (gameSize6: Est=-0.64, CrI=[-1.05, -0.25]) and more slow to improve in the six player games (block:gameSize6: Est=-0.34, CrI=[-0.43, -0.25]). Type of channel did not have a clear effect on accuracy either overall (block:channelthin: Est=-0.07, CrI=[-0.18, 0.04]) or interacting with time (block:channelthin: Est=-0.07, CrI=[-0.18, 0.04]).

The key phenomenon of iterated reference games is that the descriptions to the target images shorten over the course of repetition. This pattern held across all experiments and conditions, as the number of words produced by speakers decreased over the course of the game. The number of words produced by speakers decreases over the course of rounds (Figure ??fig:trptych-red)). In experiment 1, the overall effect of being one block later was -3.37 (CrI=[-4.54, -2.24]) words. Speakers in larger groups said more; the effect of each additional player was 1.66 (CrI=[0.66, 2.61]) words per trial, with no clear



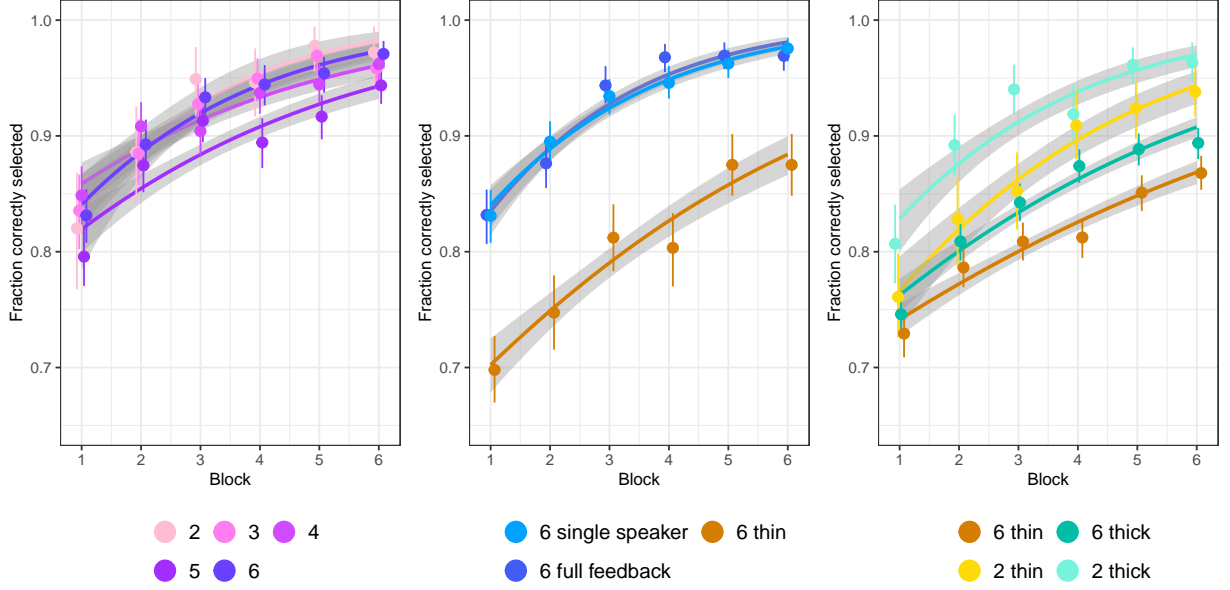


Figure 3: TODO

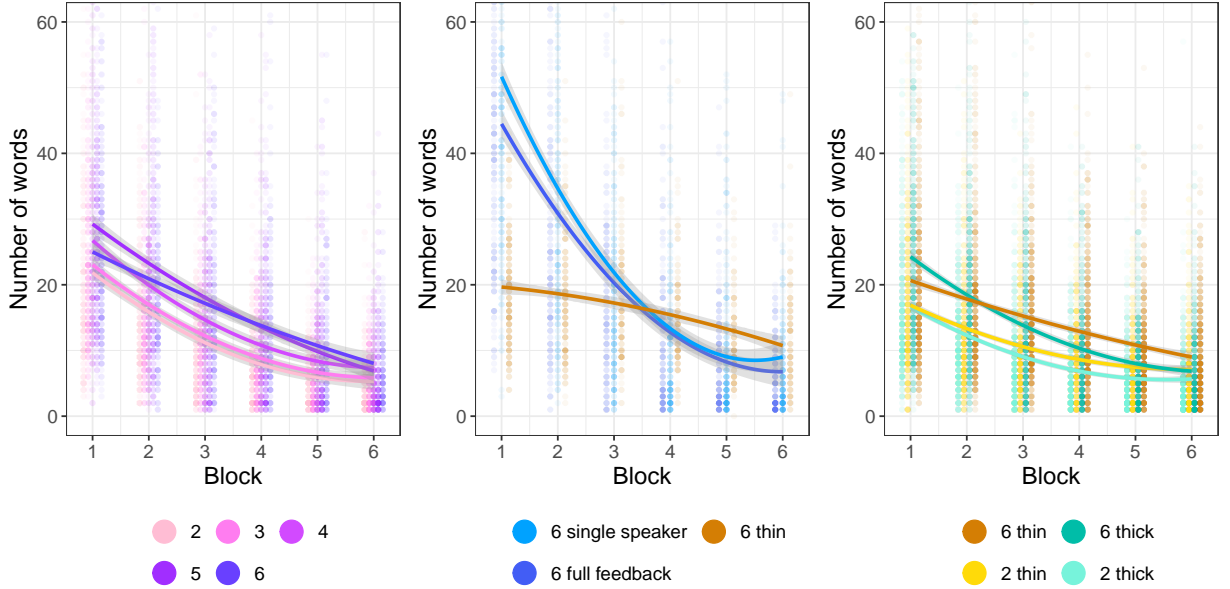


Figure 4: TODO

interaction between block and group size (block:numPlayers: Est=-0.1, CrI=[-0.36, 0.17]). In experiment 2, the result of being one block later on the number of words the speaker said per trial was Intercept: Est=29.93, CrI=[24.92, 34.84] in experiment 2a; Intercept: Est=26.03, CrI=[21.12, 30.58] in 2b, and Intercept: Est=20.5, CrI=[17.26, 23.76] in 2c. The rate of reduction was lower in the emoji condition than in the other conditions. In experiment 3, reduction occurred overall (block: Est=-2.29, CrI=[-2.95, -1.6]). The six player games said more to start with (gameSize6: Est=7.41, CrI=[3.57, 11.18]) and reduced less (block:gameSize6: Est=-1.21, CrI=[-2.06, -0.3]) than the two-player games. There were not differences due to channel type (channelthin: Est=0.63, CrI=[-3.18, 4.73]) or channel type over time (block:channelthin: Est=0.32, CrI=[-0.65, 1.24]).

TODO commentary

### 3.2 Comparisons of language between and within games

TODO describe SBERT



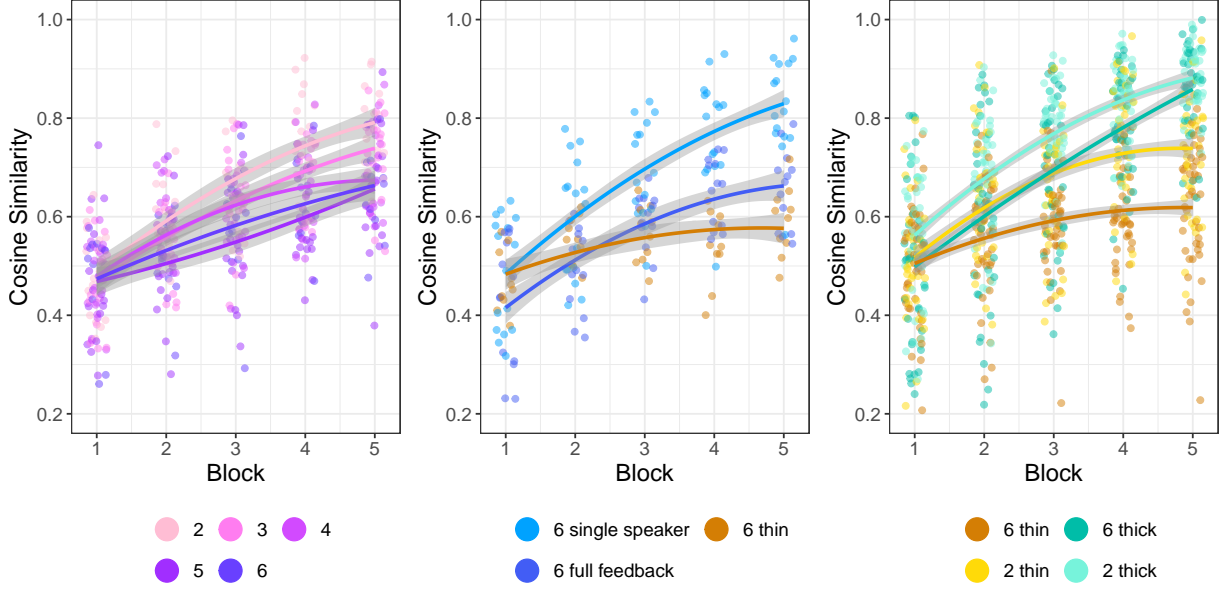


Figure 5: TODO

As a measure of convention formation, we look at the similarity of descriptions within a game for a particular tangram in different blocks. We take the last round descriptions as the established convention and measure the similarity between earlier speaker utterances and this convention.

In experiment 1, later utterances are more similar to the last utterance than earlier utterances are (0.09 (CrI=[0.08, 0.1])). The distance from the first to last utterances is invariant across group size (-0.01 (CrI=[-0.02, 0])), but smaller groups converge faster (-0.01 (CrI=[-0.01, -0.01])). Experiment 2 shows similar patterns of utterances become more similar to the last utterance, especially in the non-rotating 2a (0.09 (CrI=[0.08, 0.09])), but also in 2b (0.06 (CrI=[0.05, 0.07])) and to a smaller extent, 2c (0.02 (CrI=[0.01, 0.03])).

In terms of convergence towards the last round utterance, this is seen overall (0.08 (CrI=[0.07, 0.09])). The convergence is slower in thin games (-0.02 (CrI=[-0.03, -0.02])) and especially thin 6 player games (-0.04 (CrI=[-0.05, -0.02])).

Comparing utterances between adjacent rounds reveals similar patterns. Thin games have lower similarity between adjacent blocks (-0.12 (CrI=[-0.16, -0.09])) as do larger games (-0.03 (CrI=[-0.07, 0])). Later in the game adjacent blocks are more similar than earlier adjacent blocks (0.05 (CrI=[0.04, 0.05])), painting an overall nonlinear convergent pattern (as seen in figure TODO).

The last measure of how utterances change within games is how they compare to the first utterance; this is less good because the first utterance has more fluffy language so is less diagnostic, but later utterances are further from the first round utterance than earlier utterances (-0.03 (CrI=[-0.04, -0.03])). (TODO it's in the pre-reg, but we could dump it in a supplement? )

### 3.2.1 Divergence

Over the course of the game, different groups develop their own ways of identifying the different figures. As groups diverge, their descriptions become less similar to those used to describe the same figure in other games (-0.04 (CrI=[-0.04, -0.03])). Group size does not affect the cross-groups similarities in the first block (0 (CrI=[0, 0])), but smaller groups diverge from each other faster than larger groups (0 (CrI=[0, 0])).

All three sub-experiments show the pattern of divergence where descriptions are less similar across groups in later blocks than in earlier blocks. This reduction is stronger in 2a (-0.04 (CrI=[-0.04, -0.04])) and 2b (-0.04 (CrI=[-0.04, -0.04])) than in 2c (0 (CrI=[-0.01, 0])).

In a model of similarities between utterances produced at the same point in the game for the same tangram across different games ( $\text{sim} \sim \text{block} \times \text{channel} \times \text{gameSize} + (1|\text{tangram})$ ), games get less similar over

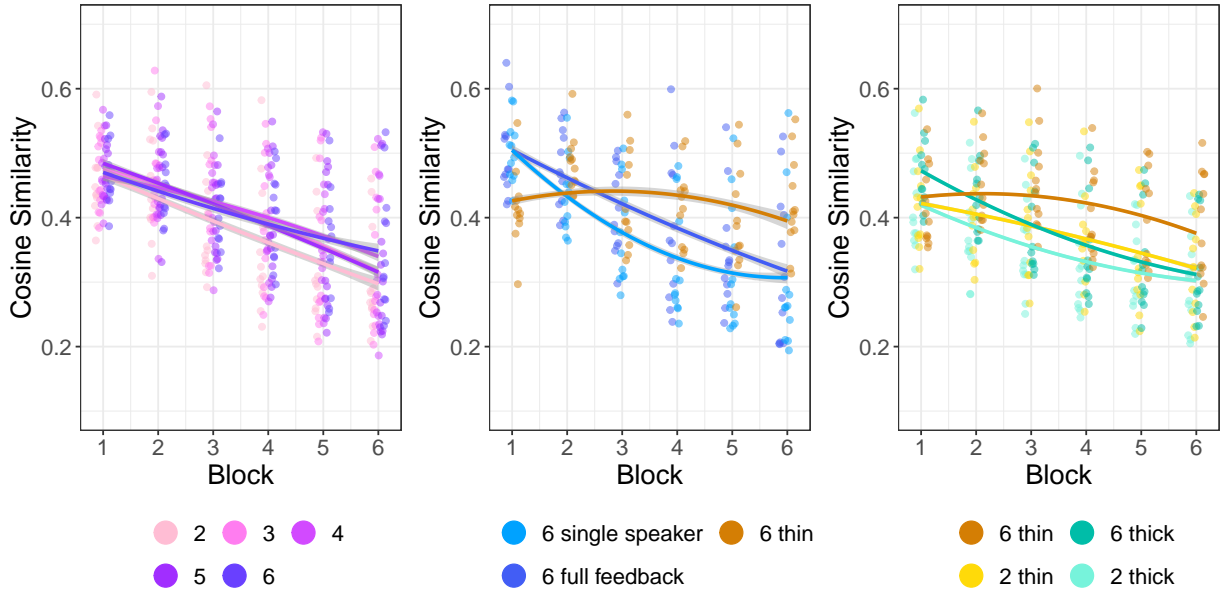


Figure 6: TODO

time (-0.02 (CrI=[-0.02, -0.02])). There are slight differences in the initial starting points across the different conditions, as well as slight condition differences in how fast the games diverge. In particular, 6 player thin games diverge more slowly (`r stats_text(div_3, 3)`).

### 3.2.2 Distinctiveness of tangrams

Another way of looking at how language changes over the course of the game is looking at how games start to refer to different tangrams more differently. This could reflect initial overlap in describing many figures as sitting or standing or by leg and arm and head position.

Over the course of the game, descriptions for each tangram become more distinctive (-0.04 (CrI=[-0.05, -0.04])). In all three subexperiments, the descriptions of tangrams become more distinctive within games across time. (2a -0.05 (CrI=[-0.05, -0.04]), 2b -0.03 (CrI=[-0.03, -0.02]), 2c -0.03 (CrI=[-0.03, -0.02])).

Tangram distinctiveness within games increased over time (-0.03 (CrI=[-0.03, -0.03])). There might be more to say about other effects, but it's mostly a starting places being different in larger games and then the slopes also differ a bit?

### 3.2.3 Emoji usage ?

TODO there's a note saying to rerun these models for longer with more extensive mixed effects!! (at least of reduction model)

### 3.2.4 Examples of names

While most groups did form conventions for most tangrams, it's illustrative to look at a case where a group did not. Table 3 shows the transcript of a 4-person group for a specific figure where they described it geometrically every round, leading to long and not very informative descriptions. Nearly all the figures have diamond heads, so this isn't a distinguishing feature, yet it is described. This illustrates the variability between groups, but also why conventions might be useful.

Table 3: Excerpt from a group that did not reduce very much. The speaker for each round is marked with (S). Figure under discussion is row 3, column 3 in Figure 1A.

Block	Person	Text
1	A(S)	Diamond on top. Body with no real arms or legs. The body is shaped like a boot with the diamond on top.
2	C B(S)	Is the boot pointed left or right? diamond on top, large body beneath it. Left is a straight line all the way down, small variations on the right to the main body
3	C(S)	Diamond in center on top. Left side straight, right side carved out like a vase.
4	D(S)	Diamond head, flat topped body, straight on the left side with two triangles pointing out on the left
5	D(S) A(S)	*on the right Diamond on top. Left side is straight, right side is obstructed, looks like a boot
	B A(S)	what do you mean by obstructed? The left side of the body is right, right side has bents in it
6	B(S)	Diamond on top of a long large body/rectangle. Left side is complete, right side has bits missing

A different 4-person group had a member who during the first block shared the idea that the task would be easier if they explicitly gave “codenames” to the figures. The transcript for this group and one of the tangrams is shown in Table ???. Of note, multiple speakers forget the assigned codename, demonstrating that meta-knowledge doesn’t always help. This group also describes the figure in relation to another already-named figured. Nonetheless, the group successfully conventionalizes on a couple reduced names for this figure: “zigzag” and “beggar”. This dual-naming of figures from multiple conceptual angles contributed by different speakers also occurs in other games.

## 4 General Discussion

this isn’t the only group dynamic; could imagine situations where listeners can see each others work collaborate (point to each other what they think, perhaps see feedback from speaker to one listener) which might make things reduce much faster

The emergence of conventions has been a key case study for communication more broadly. Yet this issue has – for the most part – been studied only in dyadic communication. While some studies have examined aspects of convention formation in larger groups (e.g., Yoon & Brown-Schmidt 2014, Yoon & Brown-Schmidt 2019), basic descriptive work has not yet investigated how group size changes the dynamics of interaction in a standard referential communication task, in part because such tasks can be difficult to administer to larger groups. Taking advantage of a new online multi-player experiment platform, we ran repeated reference games with groups of 2–6 players and characterized the nature of group performance.

Consistent with dyadic games, listeners’ selection accuracy increased over blocks at the same time as listeners sped up their selections (question 1). Crucially, speakers reduced the length of their descriptive utterances as they conventionalized on concepts for each image (question 2). Because speakers rotated, this reduction finding is robust: not only did speakers say less in later repetitions than they themselves said earlier, speakers later in the order said less than speakers earlier in the rotation. This reduction varied

with group size; smaller groups used shorter utterances, but group size did not significantly interact with block (question 3). The trajectory of reduction also depended on whether the current speaker correctly identified the tangram in the prior block and whether the current speaker was new to being speaker. This pattern is consistent with both the ‘aim low’ and ‘aim middle’ hypotheses from previous work (Yoon & Brown-Schmidt 2014, Yoon & Brown-Schmidt 2019).

What was specifically different across group sizes? Smaller groups showed more agreement in how each tangram was identified across blocks (question 4), coming to consensus earlier: Their overlap between descriptions in the first 5 blocks to the final block was higher, and words in the final block tended to originate earlier. The greater diversity in how tangrams were described in larger groups could be explained by slower convergence to a convention or parallel competing conceptualizations favored by different speakers. Larger groups have more people for the speaker to communicate to, but also more people who might interrupt with questions, and more people who have opinions about what each image looks like. Bigger groups differ from smaller groups in a number of ways, however, and disentangling these differences is an area for future work.

Group interactions are rich, and this experiment is necessarily a schematic simplification with a number of limitations. Real-life situations vary widely in who the interlocutors are, their relationships, their goals, and their environment (Carletta et al. 1998, Fay et al. 2000). Our participants were a convenience sample of Prolific workers who were strangers to each other; thus we miss richness that could come from prior relationships or shared community. Reference is only one goal out of many possible communicative goals, and the tangram images are artificial. We provided less feedback than previous studies such as Hawkins et al. (2020); this regime imitates situations where interlocutors can’t show each other examples, but it’s not representative of all communicative environments. Further, our text-based online paradigm meant that participants’ individual identities were not especially salient. In sum, communication takes place in a plethora of situations; our experiment provides some insights, but also misses many complexities that should be a focus of further experiments.

The experimental paradigm presented here could be a valuable tool to disentangle the mechanisms of group size and determine which design parameters are relevant to reduction. Luckily, with an online implementation, recruiting for and running experiments is feasible, and thus it will be possible to iterate on this experiment to determine how far the patterns generalize. While much is left to be explored, this initial data set provides a rich corpus of how humans adapt language dynamically to communicate.

## 4.1 Limitations

## 5 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*
- 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*
- 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](#)
- 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*
- 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](#)
- 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](#)

- 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](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)
- 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 (2019) Audience Design in Multiparty Conversation. *Cogn Sci* **43**:e12774. doi:[10.1111/cogs.12774](https://doi.org/10.1111/cogs.12774)