

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

Group communication is ubiquitous, but it presents some challenges not found in two-person communication, which is the primary focus of the experimental literature. One test case for communication that 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 can 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, reduction and convergence to shared descriptions occurs, but there is a gradient where smaller groups and groups with thicker communication channels show faster and stronger convergence than larger games or groups with thinner communication channels.

TODOs

* Fix first sentence of abstract (somehow?) * how to prep for interaction structure in intro

- Fix tense issues everywhere!
- Fill in some missing citations
- !!!!front load take aways in results-cosine and results-not-cosine
- Reoutline takeaways, propogate to set up/intro, results, GD
- everything gets a topic sentence
- Clean up supp figs / exile supplement

Communicating in groups can be challenging. Listeners may have different levels of background knowledge that the rest of the groups may be unaware of. Some interlocuters may interrupt with questions, and others might pipe in to explain their views, collectively leading to everyone talking at once. Multiple conversations threads may split off that need to merge back together for the group to reach agreement. Different people may understand the same speaker as meaning different things, resulting in disagreements and misunderstandings. Disagreements may even escalate to the point where meta-discussion is needed to define terms or struction the conversation differently.

TODO FIX We've all been in situations that look like this, where a conversation with half-a-dozen people devolves into chaos and results in inefficient communication. Yet, despite all of these impediments, we often communicate successfully in groups. How?

*Corresponding author. Email: vboyce@stanford.edu

One key requirement for efficient communication in groups of any size is a shared vocabulary, or shared mappings between linguistic units and objects or concepts (Traum 2004, Ginzburg & Fernandez 2005, Branigan 2006). Because reference is a requirement of communication that can be isolated and tested in experimentally manipulated contexts, it has been a case study for efficient communication more broadly. In many cases, there are widely shared convention mappings between objects and descriptions that people can rely on, but in other 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 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’. Not only are later utterances shorter than earlier utterances, but later utterances are a tacitly agreed upon name, understandable within the group, but different from the conventions chosen in other groups.

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, 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.)

How generally does this reduction phenomena occur?

What aspects of conversational infrastructure are needed to support this convergence pattern? Can a speaker form a convention when they are interacting with several listeners at once? Do listeners need to be able to ask questions or get feedback on the correct answers for the convention to form? Can conventions still form when a different person is the speaker each block? Even if it is possible to form conventions under these circumstances, is it harder to, compared to a speaker-listener pair with consistent roles and a wide communication channel?

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. Our key question was how do differences in communication channel structure impact the formation of conventionalized names in an iterated reference game. We analysed the results using traditional metrics of accuracy and number of words produced, and we used computation measures of semantic similarity to understand how utterances evolve during the games.

We find that the characteristic pattern of increasing accuracy and decreasing utterance length that is noted in the dyadic literature also occurs across conditions in multi-player games. Reduction is coupled with semantic shifts as utterances converge toward the eventual conventionalized name and away from descriptions used by other groups or for other tangrams. These convention-formation phenomena emerge across disparate conditions; however, larger groups and groups with narrower communication channels are less effective and converge more slowly and weakly than other groups.

Claims

How does the formation of shared conventions proceed in multi-party communication?

Behaviorally, larger games are similar to smaller games, but their speakers tend to say more overall, perhaps related to the increased number of listeners to respond to.

However, the interaction structure of a group affects how rapidly groups converge to partner-specific conventions, whether measured via reduction or via semantic convergence. We can interpret the experiment 3 pattern as saying that 2-player games can cope with limited feedback mechanisms, but 6-player games suffer without access to more feedback.

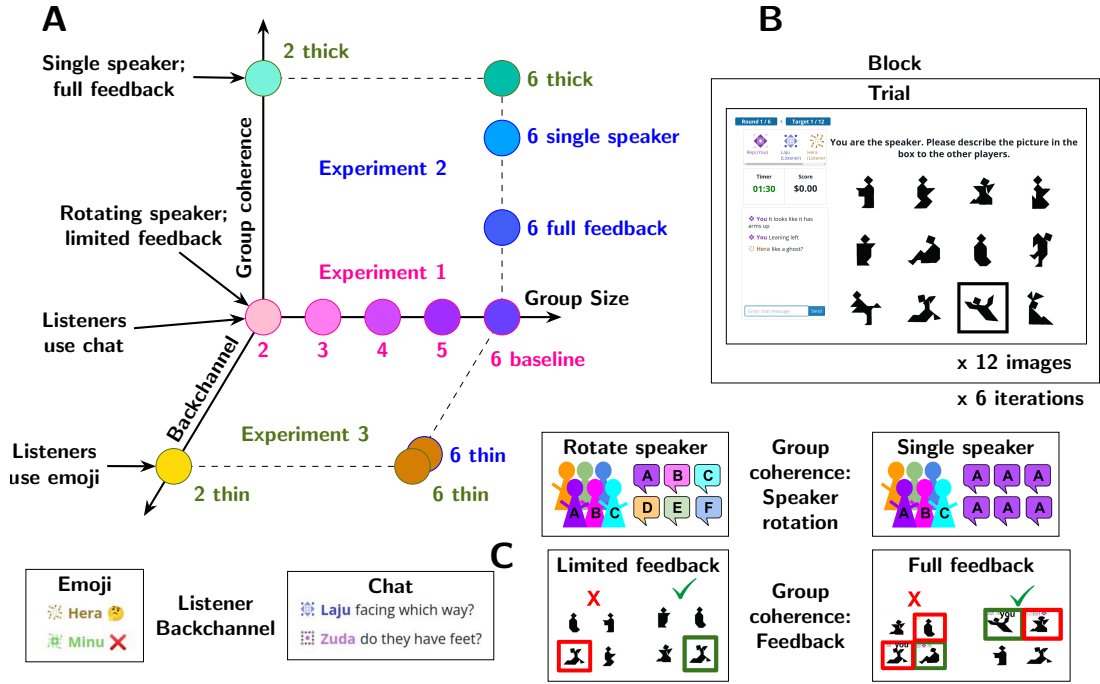


Figure 1: TODO: needs some visual help!. A: Diagram of the experimental space explored in these experiments. Experiments varied along 3 dimensions: Group size, group coherence, and listener backchannel. Each condition is shown as a dot. Experiment 1 (pink labels) varied group size from 2-6 players while holding group coherence and backchannel constant. Experiment 2 (blue labels) keep group size constant at 6 and varied the other dimensions. Relative to experiment1, 6 single 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.

Do the phenomena observed in dyadic reference games extend to larger games and games with more constrained interaction structure? How does interaction structure and game size modulate the strength of effects? Are the effects weaker in groups with more constrained interaction structure? something something linguistic analyses – do they track the verbal patterns of similarity? Are the effects weaker in some groups than others?

Listeners accuracy rose over repetitions in all games, and approached ceiling in most conditions; however, 6 player thin games were less accurate than other games (Figure 2A).

All conditions showed reduction, where speakers produced fewer words of description in later repetitions than in earlier repetitions. Speakers with more listeners talked more across blocks, and in some cases, showed sharper reduction from their wordier starting points. This is likely related to having more listeners to respond to.

According to the traditional metrics of iterated reference games, larger games are similar to small games, except that speakers say more in larger games. Listeners do not talk as much as speakers, but they talk much less in later blocks than in earlier blocks. When there are more listeners, it is more likely that some listener will talk.

Over repetitions, speaker descriptions became increasingly similar to the final description for the target image. This convergence occurred faster in smaller and higher coherence groups, and was least strong for the 6-player thin condition.

Over repetitions, speaker descriptions in one group became increasingly dissimilar to descriptions from other groups. This divergence was faster in smaller and thicker conditions, and barely occurred in the 6-player thin condition.

Results

[mini-methods]

We extended on the dyadic paradigm of Hawkins et al. (2020) by parameterizing the experiments along a few dimensions while keeping other aspects of the experiment constant. As shown in Figure ?? inset, all of the games used the same 12 target images (Clark & Wilkes-Gibbs 1986, Hawkins et al. 2020). 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 of 72 trials, where each image was described 6 times over the course of the game.

Figure 1 schematically illustrates the dimensions of variation between games and where each condition fit in the experimental space. Game size (shown on the x-axis) which varied between 2 and 6 players groups to explore the gradual effects having a larger audience. Group coherence (y-axis) was made up of two components: speaker rotation and feedback. In low group coherence games, the speaker rotated each block, while in high group coherence games, one player was the speaker for the entire game. Rotating speakers is a more stringent test of convention formation because it compares utterances from different players, but having a single speaker adds continuity that can help hold a group together. In low group coherence games, each listener only received feedback on if they were individually right or wrong in their selection; while in high group coherence games, listeners (like the speakers in all games) saw who had selected what and what the target had been, thus ensuring people saw what referent was intended and had a sense of how everyone else was doing. Listener backchannel (z-axis) varied how listeners could communicate with the group. In high backchannel games, the listeners could type text messages to the shared chat; while in low backchannel games, listeners could send 4 discrete messages (represented as emojis) to the chat. This dimension was inspired by the claimed importance of listener contributions to convention formation [CITE].

Experiment 1 varied group size with games of 2, 3, 4, 5, or 6 players all with low group coherence and high listener backchannel. Experiment 2 held group size constant at 6 while each condition deviated from experiment 1 in one aspect: 6 single speaker and 6 full feedback changed components of group coherence and 6 thin switched to a low backchannel. Finally, experiment 3 tested 4 corners of the experimental space at larger scale, with thin (low backchannel, low coherence) and thick (high backchannel, high coherence) games with either 2 or 6 players.

We compared across all of these conditions on the two behavior measures that are the common markers of reduction in the literature: listener accuracy and speaker reduction [CITE]. Additionally, we measured the similarities between speaker’s descriptions to explore how the descriptions developed over time. For analyses, we used a Bayesian multi-level regression framework with weakly regularizing priors; listener accuracy used a logistic regression, all other analyses used linear regression [CITE BRMS] [See methods for full priors, and supplement for a full list of models and model results].

Accuracy and language quantity

The canonical findings from dyadic reference games are that, over repetitions, listener accuracy remains high and increases towards a ceiling effect while the amount of referential language decreases dramatically. Listener accuracy is a measure of how successfully groups are completing the task, that is, whether speakers are succeeding in communicating the target referent to the listeners. Reduction of speaker descriptions is an indication that speaker and listeners have formed a shared conceptualization of the target, that can be distilled into a shorter and less descriptive form, while still retaining the same level of informativity to listeners. CITATIONS

In our current 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). In experiment 1, varying group size did not have a strong effect on initial accuracy (numPlayers: -0.07 [-0.2, 0.05]) or improvement rate (block:numPlayers: -0.02 [-0.05, 0.01]). In experiment 3, the 6 player games had lower initial accuracy (gameSize6: -0.64 [-1.05, -0.25]) and were slower to improve (block:gameSize6: -0.34 [-0.43, -0.25]) than the 2 player games. Thin games were not reliably worse than thick games on initial accuracy (channelthin: -0.36 [-0.78, 0.05]) or improvement rate (block:channelthin: -0.07 [-0.18, 0.04]).

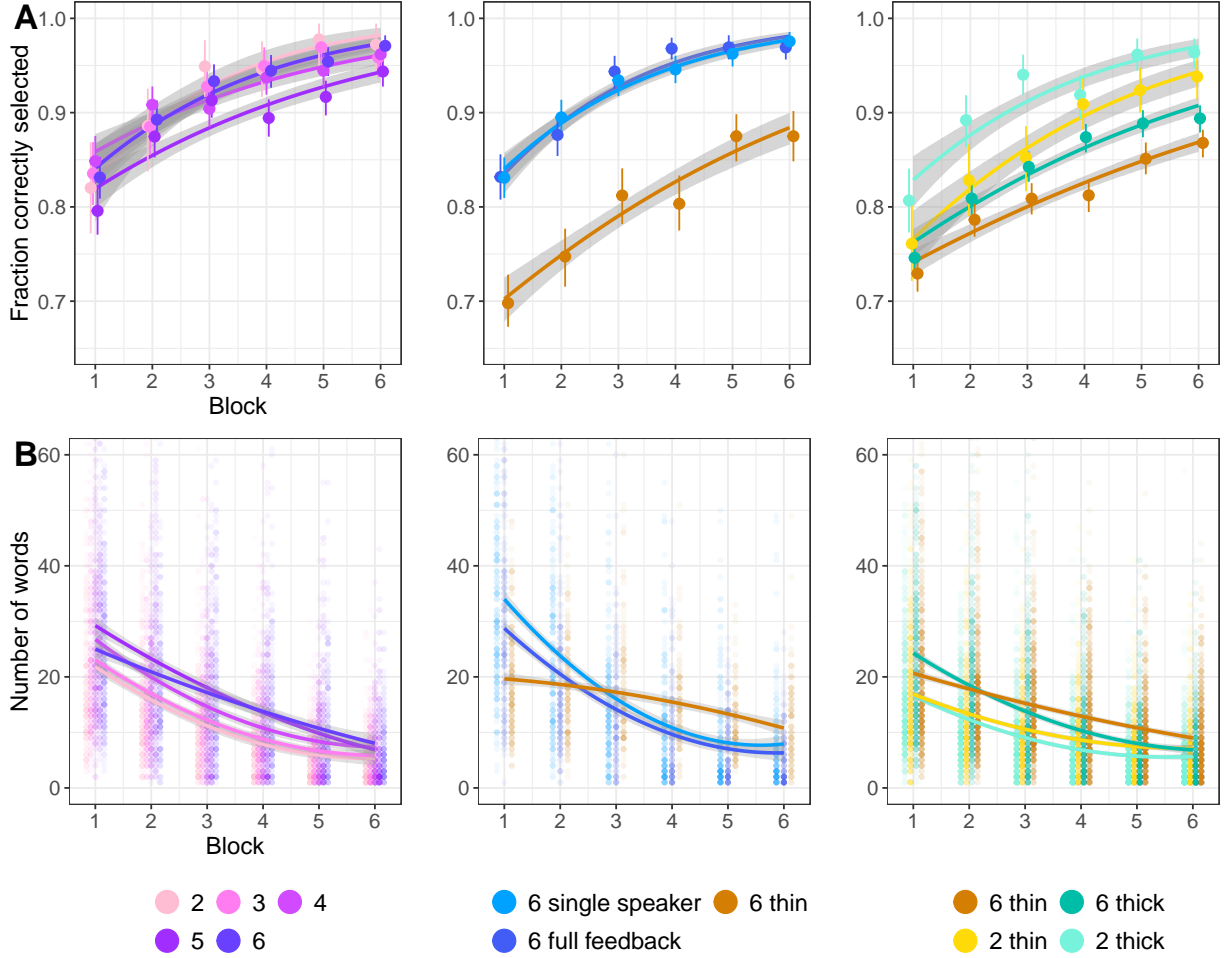


Figure 2: Behavioral results across all three experiments. A. Listener accuracy at selecting the target image. Dots are per condition, per block estimates with 95% bootstrapped CIs. Smooths are binomial fit lines. B. Number of words said by the speaker each trial. Faint dots represent individual trials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

The high and increasing levels of accuracy indicate that across all of these conditions, participants are able to succeed in communicating about the images.

The key reduction phenomena is the reduction in the number of words used to describe the images over the course of the game. The result appeared in all condition where speakers produced fewer words in later repetitions than in earlier repetitions (Figure 2B).

Speakers in larger games tended to be more verbose than speakers in smaller games, and in some cases, these speakers showed sharper reduction from initial wordiness to eventual concision. In experiment 1, the overall effect of being one block later was -3.37 $[-4.54, -2.24]$ words per trial. Speakers in larger groups said more; the effect of each additional player was 1.66 $[0.66, 2.61]$ more words per trial, with no clear interaction between block and group size. In experiment 2, the result of being one block later was block: -5.39 $[-6.46, -4.31]$ words per trial for 6 single speaker; block: -4.68 $[-5.88, -3.52]$ words per trial for 6 full feedback, and block: -2.15 $[-3.44, -1.12]$ words per trial for 6 thin. In experiment 3, the six player games said 7.41 $[3.57, 11.18]$ more words in in the first block, but reduced faster (block:gameSize6: -1.21 $[-2.06, -0.3]$) than the two-player games. Thin games were similar to thick games in initial verbosity (channelthin: 0.63 $[-3.18, 4.73]$) and reduction rate (block:channelthin: 0.32 $[-0.65, 1.24]$).

Listeners talk much less frequently than speakers, but their role in convention formation is thought to be important CITATION.

Listeners descriptions reduced over the course of the game, with more language occurring in larger games (CROSSREF SUPPLEMENT IMAGE). In experiment 1, the number of trials where any listener says

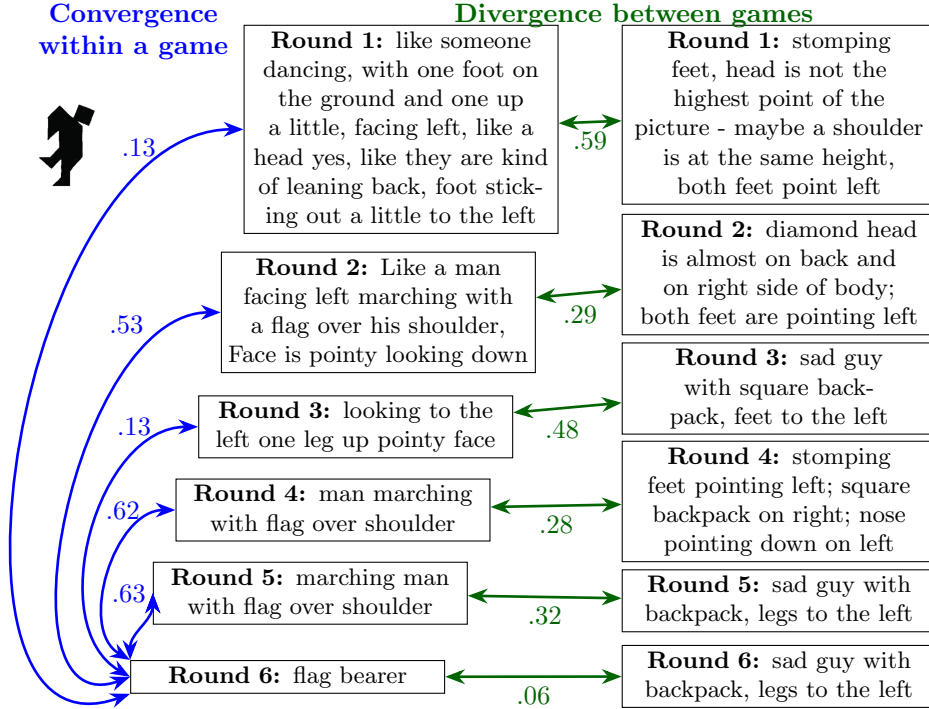


Figure 3: Example utterances describing the shown tangram figure produced by two 3-player games in Experiment 1. To measure convergence within a game (blue), we measured the cosine similarity between SBERT embeddings of descriptions and the embedding of the round 6 utterance (taken to be the convention).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

anything related to the description of the image is higher in larger groups (numPlayers: 0.78 [0.58, 0.98]), and declines across blocks (block: -0.8 [-0.97, -0.63]). When referential language is produced, more is produced in larger groups (numPlayers: 2.12 [1.03, 3.12]), but the difference in group size closes in later blocks (block:numPlayers: -0.41 [-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.

[could say more, but not sure how much to say] In games where listeners could only send 4 emoji, emoji use is common in the 6 player games, but also decreases over the course of the game (CITE FIGURE).

According to the traditional metrics of iterated reference games, larger games are similar to smaller games, except with more talking, especially early in the games. At a given level of understanding, there’s a accuracy-reduction trade-off: saying more may lead to more of the group understanding the referent, but it takes time and words. Larger groups seem to generally take the time to elaborate and increase accuracy, 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 analyses

The reduction in words is thought to be a signal of the formation of a partner-specific conceptualization of the target image that leads to an efficient way to refer to it CITATIONS. We can measure this concept formation more directly by looking at shifts in the semantic similarity of descriptions. As a group converges to a nickname for a target, descriptions within a game, to the same target should become more similar to each other and to the eventual convention. At the same time, the partner specificity means that as different groups latch onto different features as the key concept, descriptions of the same image from different groups should decrease in similarity over time.

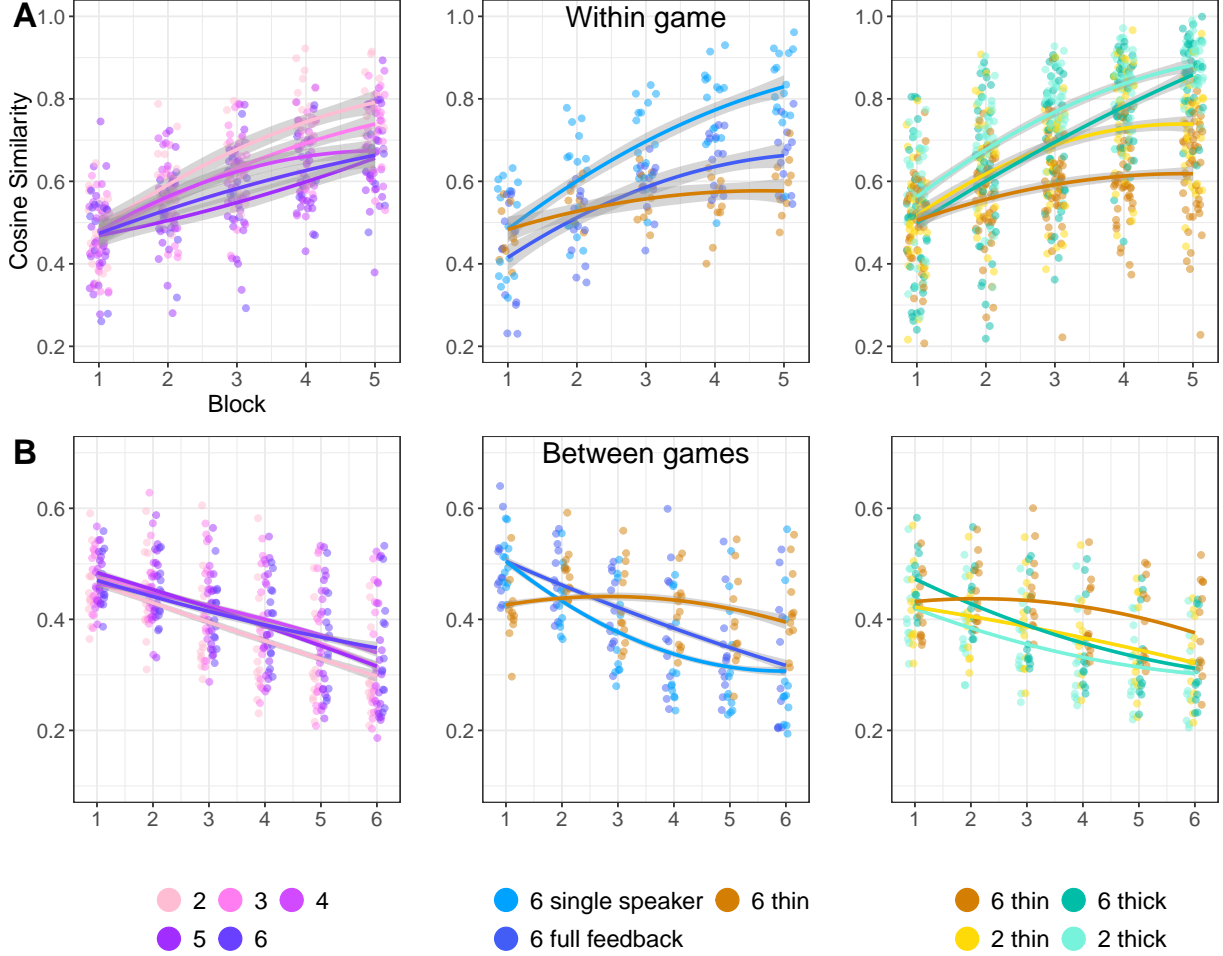


Figure 4: Language similarity results measured with pairwise cosine similarity between embeddings of two utterances. A. Convergence of utterances within games as measured by similarity between an utterance from block 1-5 to the block 6 utterance in the same game for the same image. Dots are per-game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

TO quantify the level of similarity between descriptions, we concatenated the speaker’s messages within a trial together and embedded the description into a high-dimensional vector space using SBERT (CITE). Then, we compare the similarity between any pair of utterances by taking the cosine similarity between their embeddings. These two metrics are illustrated with examples in Figure 3.

Across conditions, speaker descriptions converged toward the final description, as shown Figure 4A. This convergence occurred faster in smaller and higher coherence groups, and was least strong for the 6-player thin condition.

In experiment 1, the similarity of the first utterance to the last last utterance was invariant across group size (-0.008 [$-0.021, 0.005$]), but for smaller groups, similarity increases rapidly, and they converge faster (-0.008 [$-0.011, -0.005$]). Convergence was especially rapid in the experiment 2 single speaker condition (0.086 [$0.078, 0.094$]) where all the utterances come from the same person. In experiment 3, convergence was slower in thin games than thick games (-0.025 [$-0.033, -0.017$]) and especially thin 6 player games (-0.035 [$-0.047, -0.025$]). This convergence towards the last utterance is driven by cumulative increasing similarity between pairs of adjacent utterances (i.e. blocks 5 and 6 are closer to each other than block 1 is to 2), illustrated in TODO SEE SUPPLEMENT FIGURE.

Veronica is here in rewrite pass

Over repetitions, speaker descriptions in one group became increasingly dissimilar to descriptions from other groups. This divergence was faster in smaller and thicker conditions, and barely occurred in the

6-player thin condition. The complement to convergence within groups is divergence between groups, as different groups develop their own ways of identifying the different figures. In experiment 1, descriptions become less similar to those used to describe the same figure in other games (-0.035 [-0.038, -0.032]). Group size does not affect the cross-groups similarities in the first block (0.002 [0, 0.004]), but smaller groups diverge from each other faster than larger groups (0.001 [0.001, 0.002]). In experiment 2, divergence is stronger in the single speaker (-0.041 [-0.043, -0.039]) and full feedback conditions (-0.038 [-0.04, -0.035]) than in the 6 thin condition (-0.004 [-0.006, -0.001]). In experiment 3, descriptions from different games get less similar over time in the 2-player thick condition (-0.024 [-0.025, -0.023]). The 6 player thick condition has initially higher similarity (0.051 [0.047, 0.055]), but diverges faster (-0.008 [-0.01, -0.007]). The two-player thin has slightly higher similarity (0.014 [0.01, 0.018]) and slightly slower divergence (0.004 [0.002, 0.005]). Most noticeably, the 6 player thin games diverge much more slowly than the other conditions (0.017 [0.015, 0.019]).

As group narrow in on descriptions for each image that differ from other groups, the names for each tangram differentiate with a group (See supp figure TODO).

These semantic patterns are all consistent with the partner-specific convention development and differentiation proposed by CITE someone must have verbally described the pattern.

Discuss how the condition patterns of reduction and condition patterns of SBERT aren't quite the same!!

General Discussion

TODO fix me! Communication often occurs in multi-party settings, but much of the research on communication has been on dyads. One common paradigm for studying communication has been iterated reference games, where the phenomenon of reduction over repeated reference has been robustly established. How does the formation of shared conventions proceed in multi-party communication?

We build on this tradition by expanding the reference game paradigm to larger groups under various 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 shows the same patterns of increasing accuracy, reduction in speaker utterances, and semantic convergence within and divergence between groups. Thus, the classic reduction phenomena can occur even under situations with constrained interaction structure.

However, the interaction structure of a group affects how rapidly groups converge to partner-specific conventions, whether measured via reduction or via semantic convergence. 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. We can interpret the experiment 3 pattern as saying that 2-player games can cope with limited feedback mechanisms, but 6-player games suffer without access to more feedback.

The 6 thin condition also presents intriguing evidence that reduction and partner-specific utterance formation might not be closely tethered: this condition shows reduction that is not substantially different than in the 6 thick condition, but it shows much less semantic convergence than other conditions. This is just one condition, but it suggests that descriptions can shorten without differentiating, perhaps all shortening to a common prior. CITE THAT PAPER ABOUT GROUP SIZE This highlights the importance of measuring semantics separately from just length proxies. Our goal here was to further explore the space; not to directly adjudicate theories. However, given how little listeners talked, and the success of games where listeners had a limited backchannel, it seems that people are readily able to jointly agree on conventions without talking (consistent with RSA-etc CITATIONS models).

Overall, we have a number of groups, and a large corpus of language that could be probed further. The number of groups in any one condition is limited, but overall, there's a lot of people and a lot of language. TODO needs work

[not sure if/where to include] My biggest take away from this is that the hard part seems to be getting first contact with a shared meaning (and knowing you have done this), and the "reducing" part is easier since you have something to work with. So, how ideas of how to describe are sampled might be worth exploring more?

Limitations

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. Group to group variability is high, as individual differences in communication style and shape description ability interacted. Coupled with small numbers of groups in some conditions, measures may be somewhat noisy depending on the random assignment of participants to games and conditions.

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 differences in group size that could then lead to more experimental conditions.

Conclusion ??

something something about how set-ups like this could explore related phenomena picking other points in reference game / ref game adjacent space

TODO

Methods

For all experiments, 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 ??A) 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. 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

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: single 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>

Table 1: 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	Total Participants
1: baseline	2	15	4	38
1: baseline	3	18	2	60
1: baseline	4	19	2	84
1: baseline	5	17	3	100
1: baseline	6	12	6	108
2: single speaker	6	15	3	108
2: full feedback	6	13	4	102
2: thin	6	10	6	96
3: thin	2	35	3	76
3: thin	6*	44	0	235
3: thick	2	39	3	84
3: thick	6*	38	2	222

1319 people participated across the 3 experiments. A breakdown of number of games and participants in each condition is shown in Table 1.

Materials

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

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. 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.

Once the game started, participants saw screens like Figure ??A. 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 ??B). 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. 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 single 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 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.

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 plays in these 6* player games is at [TODO SUPPLEMENT!](#) The realities of online recruitment and disconnection meant that the number of games varied, although we aimed for 20 games in each condition in experiments 1 and 2, and 40 per condition in experiment 3. We excluded incomplete blocks from analyses, but included complete blocks from partial games (See [Table 1](#) 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

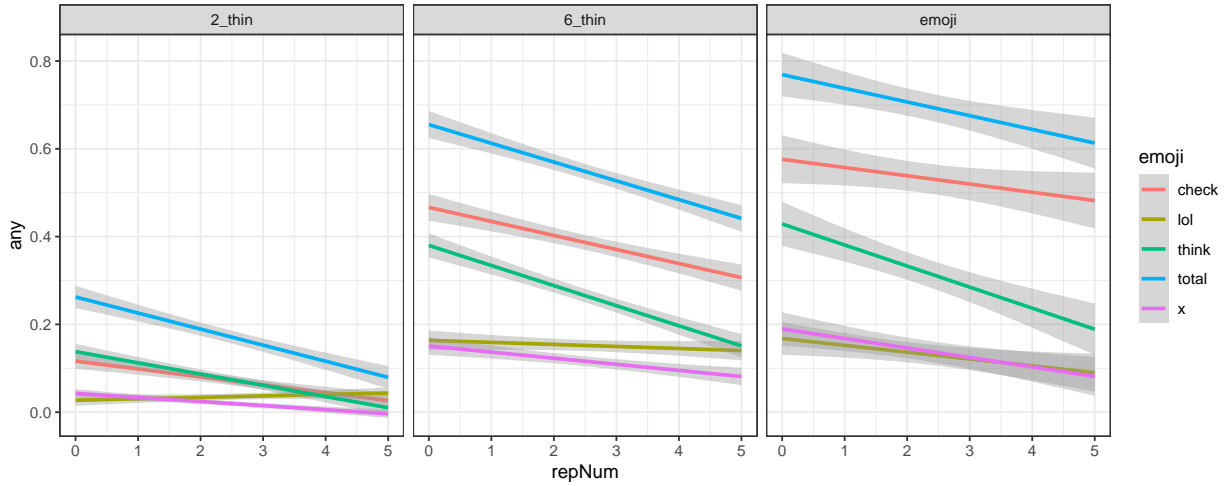


Figure 5: Fraction of trials on which at least one listener produced the labelled emoji (or any emoji).

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 (CITE) 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 results used a logistic model, other results use linear models. Accuracy models were all 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.

Supplement

This supplement will get moved to a separate file, but for now is tacked on.

Figures of listener performance

Emoji use is common in the 6 player thin games where most trials at least one listener used at least one emoji, but emoji use declines over blocks. SEE FIGURE WHATEVER IN SUPPLEMENT The use of emoji in the thin games is not directly comparable to listener contributes, since some emoji usage (such as the green checkmark) are most likely equivalent to non-referential listener language (“got it” etc.) that was excluded. The higher rate of emoji use versus referential language could be due to it’s non-equivalence, a lower level of accuracy in thin games, or emojis being a lower threshold for sending than written out questions.

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 (-0.043 [-0.046, -0.039]). In all three subexperiments, the descriptions of tangrams become more distinctive within games across time. (2a -0.046 [-0.048, -0.044], 2b -0.025 [-0.028, -0.022], 2c -0.025 [-0.028, -0.022]).

Tangram distinctiveness within games increased over time (-0.027 [-0.029, -0.025]). 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?

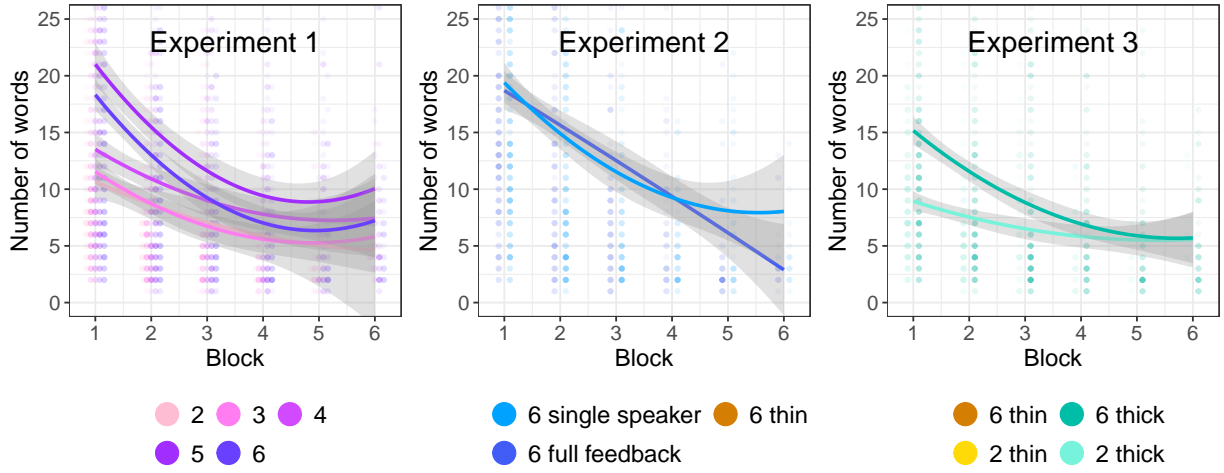


Figure 6: Number of words of referential language produced by listeners over time. Excludes trials where no listeners contributed descriptive language.

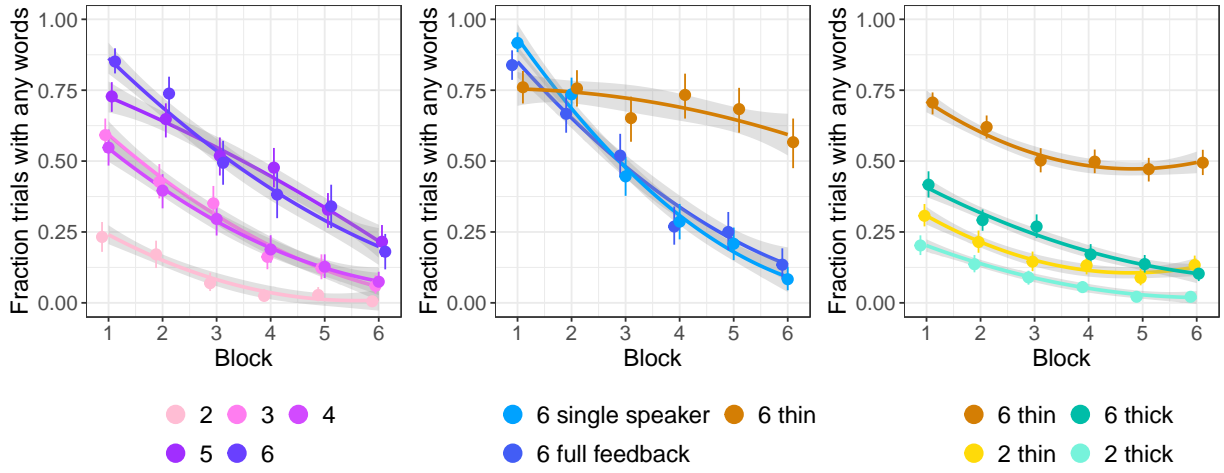


Figure 7: Fraction of trials when any reference language (or emoji) was produced by any listener.

play with more diagrams

Comparing utterances between adjacent rounds reveals similar patterns. Thin games have lower similarity between adjacent blocks (-0.124 $[-0.159, -0.088]$) as do larger games (-0.034 $[-0.069, 0.003]$). Later in the game adjacent blocks are more similar than earlier adjacent blocks (0.046 $[0.041, 0.052]$), painting an overall nonlinear convergent pattern (as seen in Figure 8).

Models

Note that for all models, block was 0 indexed, so intercepts are what happened during the first block.

Accuracy models

Accuracy models were all run as logistic models with $\text{normal}(0,1)$ priors for both betas and sd. This model was not explicitly included in the experiment 1 and 2 pre-registrations; it was included with more ambitious mixed effects (which did not run in a timely manner) in the experiment 3 pre-reg.

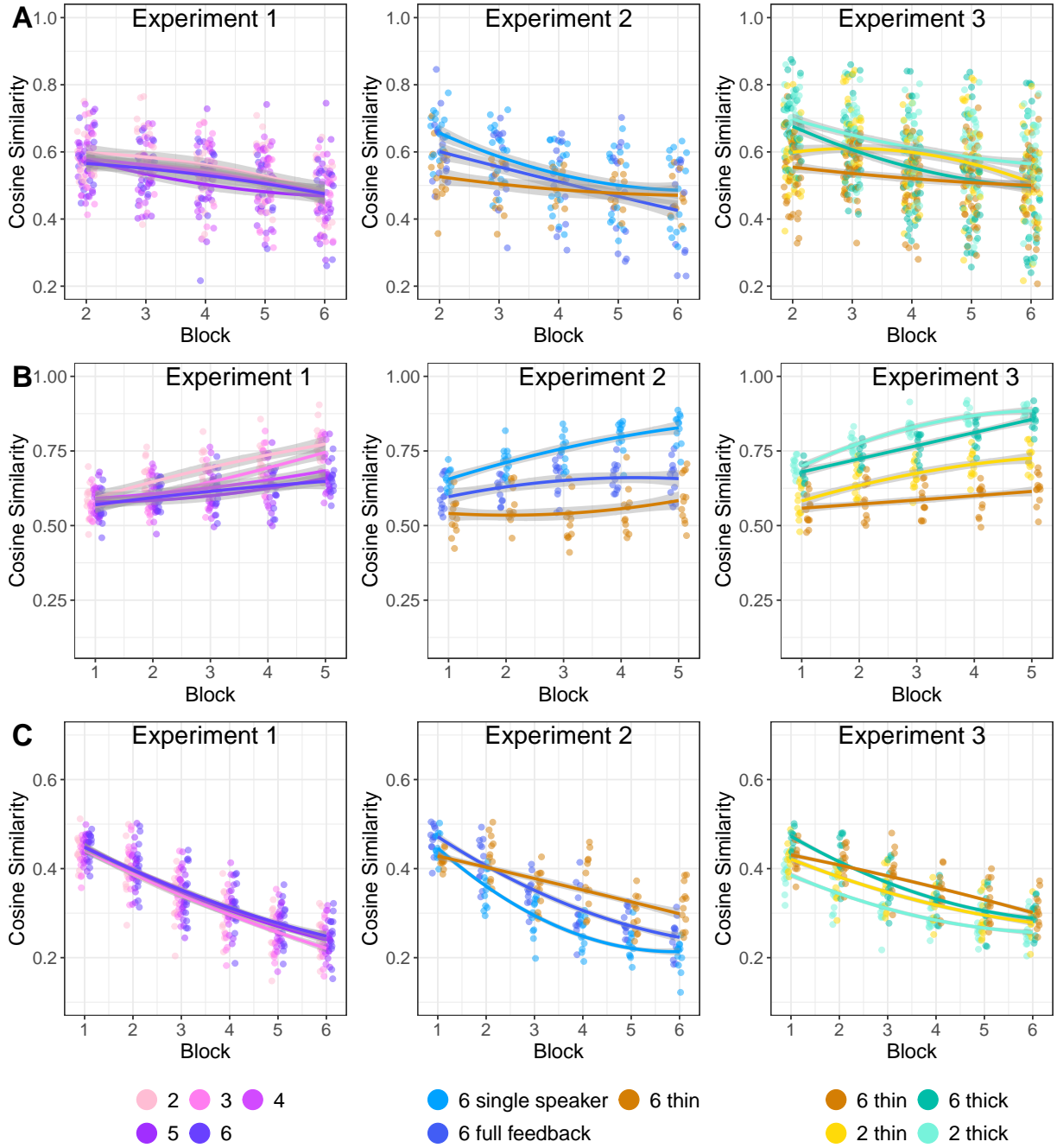


Figure 8: Additional measures of convergence and divergence. A is similarity to first utterance. B is similarity between utterances from adjacent blocks. C is divergence in descriptions of different tangrams within a group

Table 2: Experiment 1 logistic model of listener accuracy:
 $\text{correct.num} \sim \text{block} \times \text{numPlayers} + (1|\text{gameId})$

Term	Est.	CrI
block	0.44	[0.31, 0.58]
block:numPlayers	-0.02	[-0.05, 0.01]
Intercept	2.10	[1.57, 2.65]
numPlayers	-0.07	[-0.2, 0.05]

Table 3: Experiment 2: 6 single speaker logistic model of listener accuracy:
 $\text{correct.num} \sim \text{block} + (1|\text{gameId})$

Term	Est.	CrI
block	0.45	[0.39, 0.52]
Intercept	1.78	[1.4, 2.19]

Table 4: Experiment 2: 6 full feedback logistic model of listener accuracy:
 $\text{correct.num} \sim \text{block} + (1|\text{gameId})$

Term	Est.	CrI
block	0.47	[0.39, 0.54]
Intercept	1.35	[0.59, 2.06]

Table 5: Experiment 2: 6 thin logistic model of listener accuracy:
 $\text{correct.num} \sim \text{block} + (1|\text{gameId})$

Term	Est.	CrI
block	0.23	[0.19, 0.28]
Intercept	0.88	[0.64, 1.12]

Table 6: Experiment 3 logistic model of listener accuracy:
 $\text{correct.num} \sim \text{block} \times \text{gameSize} \times \text{channel} + (1|\text{gameId})$

Term	Est.	CrI
block	0.41	[0.32, 0.5]
block:channelthin	-0.07	[-0.18, 0.04]
block:gameSize6	-0.34	[-0.43, -0.25]
block:gameSize6:channelthin	0.07	[-0.05, 0.19]
channelthin	-0.36	[-0.78, 0.05]
gameSize6	-0.64	[-1.05, -0.25]
gameSize6:channelthin	0.31	[-0.22, 0.87]
Intercept	1.69	[1.39, 1.99]

Reduction models

Reduction models were run as linear models with an intercept prior of $\text{normal}(12,20)$, a beta prior of $\text{normal}(0,10)$, an sd prior of $\text{normal}(0,5)$ and a correlation prior of $\text{lkj}(1)$. This model was pre-registered for each experiment and run with the mixed effects structure as prespecified.

Table 7: Experiment 1:
 $\text{words} \sim \text{block} \times \text{numPlayers} + (\text{block}|\text{tangram}) + (1|\text{playerId}) + (1|\text{tangram:gameId}) + (\text{block}|\text{gameId})$

Term	Est.	CrI
block	-3.37	[-4.54, -2.24]
block:numPlayers	-0.10	[-0.36, 0.17]
Intercept	16.79	[11.96, 21.93]
numPlayers	1.66	[0.66, 2.61]

Table 8: Experiment 2: 6 single speaker:

words \sim block + (block|tangram) + (1|tangram:gameId) + (block|gameId)

Term	Est.	CrI
block	-5.39	[-6.46, -4.31]
Intercept	29.93	[24.92, 34.84]

Table 9: Experiment 2: 6 full feedback:

words \sim block + (block|tangram) + (1|tangram:gameId) + (block|gameId)

Term	Est.	CrI
block	-4.68	[-5.88, -3.52]
Intercept	26.03	[21.12, 30.58]

Table 10: Experiment 2: 6 thin:

words \sim block + (block|tangram) + (1|tangram:gameId) + (block|gameId)

Term	Est.	CrI
block	-2.15	[-3.44, -1.12]
Intercept	20.50	[17.26, 23.76]

Table 11: Experiment 3:

words \sim block \times channel \times gameSize + (block \times channel \times gameSize|tangram) + (1|tangram:gameId) + (block|gameId)

Term	Est.	CrI
block	-2.29	[-2.95, -1.6]
block:channelthin	0.32	[-0.65, 1.24]
block:channelthin:gameSize6	0.64	[-0.61, 1.89]
block:gameSize6	-1.21	[-2.06, -0.3]
channelthin	0.63	[-3.18, 4.73]
channelthin:gameSize6	-2.11	[-7.41, 2.98]
gameSize6	7.41	[3.57, 11.18]
Intercept	14.99	[11.86, 17.89]

Extra reduction model

For experiment 1, we also pre-specified models about whether the speaker’s correctness (as a listener) on the prior block had an effect

Model of whether speaker’s correct/incorrect answer in previous block has an effect

```
## [1] "words~$\\sim$ block~$\\times$ numPlayers~+ block~$\\times$ wasINcorrect~+ (block|tangram)~+
```

Listener models

Table 12: Experiment 1:

words \sim block \times numPlayers + block \times wasINcorrect + (block|tangram) + (1|playerId) + (1|tangram:gameId) + (block|gameId)

Term	Est.	CrI
block	-2.15	[-3.35, -0.98]
block:numPlayers	-0.23	[-0.51, 0.06]
block:wasINcorrect	0.25	[-0.24, 0.74]
Intercept	11.98	[6.31, 17.7]
numPlayers	2.15	[0.93, 3.36]
wasINcorrect	3.08	[1.69, 4.42]

Table 13: Experiment 1: words \sim block \times numPlayers + (block|gameId)

Term	Est.	CrI
block	-0.17	[-1.63, 1.24]
block:numPlayers	-0.41	[-0.72, -0.09]
Intercept	4.67	[0.09, 9.32]
numPlayers	2.12	[1.03, 3.12]

Table 14: Experiment 1: is.words \sim block \times numPlayers + (1|gameId)

Term	Est.	CrI
block	-0.80	[-0.97, -0.63]
block:numPlayers	0.03	[0, 0.07]
Intercept	-2.65	[-3.5, -1.83]
numPlayers	0.78	[0.58, 0.98]

SBERT models

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.

These models were verbally described (but not formally specified) in the pre-registrations for experiment 2 in the full feedback and thin conditions and for experiment 3, for looking at divergence between games, convergence within games (compare to first, next, and last), and divergence between tangrams within games.

Convergence within games: comparison to last round

This is the convergence metric presented in the paper.

Table 15: Experiment 1: sim \sim earlier \times condition + (1|tangram) + (1|gameId)

Term	Est.	CrI
condition	-0.008	[-0.021, 0.005]
earlier	0.089	[0.076, 0.102]
earlier:condition	-0.008	[-0.011, -0.005]
Intercept	0.517	[0.458, 0.573]

Table 16: Experiment 2: 6 single speaker:sim \sim earlier + (1|tangram) + (1|gameId)

Term	Est.	CrI
earlier	0.086	[0.078, 0.094]
Intercept	0.499	[0.444, 0.556]

Table 17: Experiment 2: 6 full feedback:sim \sim earlier + (1|tangram) + (1|gameId)

Term	Est.	CrI
earlier	0.062	[0.051, 0.072]
Intercept	0.438	[0.389, 0.487]

Table 18: Experiment 2: 6 thin:sim \sim earlier + (1|tangram) + (1|gameId)

Term	Est.	CrI
earlier	0.023	[0.013, 0.033]
Intercept	0.498	[0.453, 0.54]

Table 19: Experiment 3:sim \sim earlier \times channel \times gameSize + (1|tangram) + (1|gameId)

Term	Est.	CrI
channelthin	-0.034	[-0.08, 0.011]
channelthin:gameSize6	0.039	[-0.021, 0.097]
earlier	0.080	[0.074, 0.086]
earlier:channelthin	-0.025	[-0.033, -0.017]
earlier:channelthin:gameSize6	-0.035	[-0.047, -0.025]
earlier:gameSize6	0.009	[0.001, 0.017]
gameSize6	-0.069	[-0.113, -0.025]
Intercept	0.581	[0.542, 0.62]

Divergence across games

To look at how games diverged from each other ... TODO

Table 20: Experiment 1:sim \sim block \times condition + (1|tangram)

Term	Est.	CrI
block	-0.035	[-0.038, -0.032]
block:condition	0.001	[0.001, 0.002]
condition	0.002	[0, 0.004]
Intercept	0.468	[0.429, 0.507]

Table 21: Experiment 2: 6 single speaker:sim \sim block + (1|tangram)

Term	Est.	CrI
block	-0.041	[-0.043, -0.039]
Intercept	0.484	[0.442, 0.526]

Table 22: Experiment 2: 6 full feedback:sim \sim block + (1|tangram)

Term	Est.	CrI
block	-0.038	[-0.04, -0.035]
Intercept	0.502	[0.46, 0.546]

Table 23: Experiment 2: 6 thin:sim \sim block + (1|tangram)

Term	Est.	CrI
block	-0.004	[-0.006, -0.001]
Intercept	0.434	[0.406, 0.465]

Table 24: Experiment 3:sim \sim block \times channel \times gameSize + (1|tangram)

Term	Est.	CrI
block	-0.024	[-0.025, -0.023]
block:channelthin	0.004	[0.002, 0.005]
block:channelthin:gameSize6	0.017	[0.015, 0.019]
block:gameSize6	-0.008	[-0.01, -0.007]
channelthin	0.014	[0.01, 0.018]
channelthin:gameSize6	-0.030	[-0.035, -0.024]
gameSize6	0.051	[0.047, 0.055]
Intercept	0.411	[0.368, 0.453]

Divergence across tangrams

Table 25: Experiment 1:sim \sim block \times condition + (1|gameId)

Term	Est.	CrI
block	-0.043	[-0.046, -0.039]
block:condition	0.000	[-0.001, 0.001]
condition	0.003	[-0.008, 0.014]
Intercept	0.429	[0.382, 0.473]

Table 26: Experiment 2: 6 single speaker:sim \sim block + (1|gameId)

Term	Est.	CrI
block	-0.046	[-0.048, -0.044]
Intercept	0.416	[0.389, 0.443]

Table 27: Experiment 2: 6 full feedback:sim \sim block + (1|gameId)

Term	Est.	CrI
block	-0.047	[-0.049, -0.044]
Intercept	0.459	[0.422, 0.496]

Table 28: Experiment 2: 6 thin:sim \sim block + (1|gameId)

Term	Est.	CrI
block	-0.025	[-0.028, -0.022]
Intercept	0.432	[0.393, 0.471]

Table 29: Experiment 3:sim \sim block \times channel \times gameSize + (1|gameId)

Term	Est.	CrI
block	-0.027	[-0.029, -0.025]
block:channelthin	-0.001	[-0.003, 0.002]
block:channelthin:gameSize6	0.011	[0.008, 0.015]
block:gameSize6	-0.010	[-0.013, -0.008]
channelthin	0.038	[-0.001, 0.082]
channelthin:gameSize6	-0.053	[-0.115, 0]
gameSize6	0.073	[0.035, 0.113]
Intercept	0.378	[0.352, 0.404]

convergence to next

We also looked at how similar an utterance was to the next round utterance: this can be thought of as the derivative of the to-last comparison. (although cosine similarities are not actually additive in the same way integrals are)

Table 30: Experiment 1:sim \sim earlier \times condition + (1|tangram) + (1|gameId)

Term	Est.	CrI
condition	-0.004	[-0.014, 0.006]
earlier	0.063	[0.051, 0.075]
earlier:condition	-0.008	[-0.011, -0.006]
Intercept	0.591	[0.541, 0.641]

Table 31: Experiment 2: 6 single speaker:sim \sim earlier + (1|tangram) + (1|gameId)

Term	Est.	CrI
earlier	0.043	[0.037, 0.05]
Intercept	0.660	[0.619, 0.702]

Table 32: Experiment 2: 6 full feedback:sim \sim earlier + (1|tangram) + (1|gameId)

Term	Est.	CrI
earlier	0.015	[0.006, 0.024]
Intercept	0.605	[0.569, 0.643]

Table 33: Experiment 2: 6 thin:sim \sim earlier + (1|tangram) + (1|gameId)

Term	Est.	CrI
earlier	0.010	[0, 0.019]
Intercept	0.533	[0.49, 0.578]

Table 34: Experiment 3:sim \sim earlier \times channel \times gameSize + (1|tangram) + (1|gameId)

Term	Est.	CrI
channelthin	-0.124	[-0.159, -0.088]
channelthin:gameSize6	0.000	[-0.051, 0.049]
earlier	0.046	[0.041, 0.052]
earlier:channelthin	-0.010	[-0.018, -0.002]
earlier:channelthin:gameSize6	-0.018	[-0.029, -0.007]
earlier:gameSize6	-0.003	[-0.011, 0.004]
gameSize6	-0.034	[-0.069, 0.003]
Intercept	0.714	[0.682, 0.746]

divergence from first

We also looked at how similar an utterance was to the first round utterance. This is not very informative because first round utterances tend to be pretty unwieldy. TODO explain more or don't include

Table 35: Experiment 1:sim \sim later \times condition + (1|tangram) + (1|gameId)

Term	Est.	CrI
condition	-0.010	[-0.022, 0.003]
Intercept	0.647	[0.591, 0.705]
later	-0.030	[-0.041, -0.019]
later:condition	0.001	[-0.002, 0.004]

Table 36: Experiment 2: 6 single speaker:sim \sim later + (1|tangram) + (1|gameId)

Term	Est.	CrI
Intercept	0.680	[0.628, 0.728]
later	-0.042	[-0.049, -0.035]

Table 37: Experiment 2: 6 full feedback:sim \sim later + (1|tangram) + (1|gameId)

Term	Est.	CrI
Intercept	0.644	[0.584, 0.706]
later	-0.044	[-0.052, -0.037]

Extra emoji analysis

Written about 6thin in experiment 2 and for 2 and 6 thin in 3 Additionally, exclusive to this condition, we will analyse the distribution of emoji's produced as a function of block and its relation to accuracy and speaker utterance length.

Table 38: Experiment 2: 6 thin:sim \sim later + (1|tangram) + (1|gameId)

Term	Est.	CrI
Intercept	0.537	[0.49, 0.584]
later	-0.014	[-0.023, -0.004]

Table 39: Experiment 3:sim \sim later \times channel \times gameSize + (1|tangram) + (1|gameId)

Term	Est.	CrI
channelthin	-0.076	[-0.123, -0.026]
channelthin:gameSize6	-0.062	[-0.127, 0.001]
gameSize6	-0.017	[-0.062, 0.03]
Intercept	0.721	[0.681, 0.76]
later	-0.034	[-0.039, -0.028]
later:channelthin	0.011	[0.003, 0.019]
later:channelthin:gameSize6	0.021	[0.01, 0.032]
later:gameSize6	-0.011	[-0.019, -0.004]

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*
- Branigan H (2006) Perspectives on multi-party dialogue. *Research on Language and Computation* **4**:153–177
- Clark HH, Wilkes-Gibbs D (1986) [Referring as a collaborative process](#). *Cognition*
- 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
- 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](#)
- Schober MF, Clark HH (1989) Understanding by addressees and overhearers. *Cognitive Psychology* **21**:211–232. doi:[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](#)
- 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](#)