

One-to-Many Communication and Compositionality in Emergent Communication

Heeyoung Lee

Sungkyunkwan University
Suwon, South Korea
hy18284@g.skku.edu

Abstract

Compositional languages leverage rules that derive meaning from combinations of simpler constituents. This property is considered to be the hallmark of human language as it enables the ability to express novel concepts and ease of learning. As such, numerous studies in the emergent communication field explore the prerequisite conditions for emergence of compositionality. Most of these studies set out one-to-one communication environment wherein a speaker interacts with a single listener during a single round of communication game. However, real-world communications often involve multiple listeners; their interests may vary and they may even need to coordinate among themselves to be successful at a given task. This work investigates the effects of one-to-many communication environment on emergent languages where a single speaker broadcasts its message to multiple listeners to cooperatively solve a task. We observe that simply broadcasting the speaker’s message to multiple listeners does not induce more compositional languages. We then find and analyze two axes of environmental pressures that facilitate emergence of compositionality: listeners of *different interests* and *coordination* among listeners.

1 Introduction

The field of emergent communication studies the core environmental factors in language emergence and the characteristics of emergent languages in relation to those of humans. The recent developments in artificial neural networks have spurred research on the field utilizing communication simulations of neural agents (Lazaridou and Baroni, 2020). This has served as a crucial testbed for studying evolution of language (Briscoe, 2002), which often lacks concrete physical trace. The field has also demonstrated promising application possibilities in numerous domains leveraging language’s desirable properties (Mu et al., 2023; Yao et al., 2022; Xu et al., 2022).

Compositionality (Janssen and Partee, 1997) is one of the most prominent features of human languages. Compositional languages can express complex meaning with combinations of simpler attributes leveraging systematic rule structures. This enables the ability to express novel concepts by combining familiar attributes. Compositionality is also attributed to enhancing languages’ learnability (Ren et al., 2020; Davidson, 1965) and gives rise to robustness to noisy communication channel (Kuciński et al., 2021).

Determining the prerequisite environmental pressures for emergence of compositionality has been extensively studied in the field. These factors include language’s learnability (Ren et al., 2020; Chaabouni et al., 2020; Smith et al., 2003; Li and Bowling, 2019), agents’ capacity (Resnick et al., 2020), reliability of communication channel (Kuciński et al., 2021), task difficulty (Chaabouni et al., 2022; Choi et al., 2018; Mu and Goodman, 2021; Bouchacourt and Baroni, 2018; Lazaridou et al., 2017), and communication channel capacity (Lazaridou et al., 2018; Chaabouni et al., 2020). Recently, populations of agents have been investigated as a driving force for emergence of compositionality (Rita et al., 2022a; Michel et al., 2023) following prior sociolinguistic findings that larger population sizes tend to derive more structured languages (Raviv et al., 2019).

Most of these studies take one-to-one communication regime where only a single speaker-listener pair interacts with each other during an instance of game play. Even when there are multiple listeners in the system, a speaker’s message is only sent to a single listener (Chaabouni et al., 2022; Michel et al., 2023; Rita et al., 2022a; Kim and Oh, 2021; Tieleman et al., 2019). Consequently, they fail to model the effects of one-to-many communication in emergent languages.

This work investigates the effects of one-to-many communication regime on the compositional-

ity of emergent languages. In real-world communications, a single message often concerns multiple parties: an advertisement of a product, a sergeant's command to a squad, etc. In these scenarios, there are more than one interested entity for a given message. This environment opens two interesting aspects of communication, and we find that these aspects each introduce a new environmental pressure that facilitates emergence of compositionality.

First, the listeners may not share the same interests. In the case of the advertisement of a product, some of the viewers of the advertisement may only be interested in certain characteristics of the product such as colors and sizes, while others may only care about the price and brand name. While it is still the case that the advertisement must contain all of the relevant information for the product, we argue that it introduces a new pressure that forces the message to be easier to understand for listeners that are only interested in certain parts of the attributes. We hypothesize that these listeners would prefer messages that are easily interpretable, without the need to understand other details corresponding to attributes that they are not concerned with.

Second, listeners may need to coordinate among themselves to be successful at the task at hand. In the case of the sergeant's command to a squad, coordination among the squad may be required for them to have successfully carried out the mission. Hence, a misinterpretation of the command from a single listener may result in failure for the entire squad. We argue that the pressure that the language be simultaneously understood by multiple listeners forces the language to be more compositional. Intuitively, it is plausible that one listener may develop a compositionally inferior language, but it is less likely to be shared by other listeners in the group due to its inferior compositionality.

Extensive experiments confirm the hypotheses that agents of different interests and coordination among agents are crucial environmental pressures for emergence of compositionality. We find that simply broadcasting a speaker's message to multiple listeners does not enhance compositionality of induced languages. We observe emergence of compositionality when listeners of different interests are introduced or coordination pressures are injected to the environment. We then analyze what kinds of compositionality structures are derived from these pressures with various compositionality measures.

2 Related work

Emergent communication and its applications

Human languages exhibit a number of universal characteristics (Greenberg, 1961). The emergent communication field strives to close the gap between the communication protocols emerged from artificial agents and the natural languages with regard to these language universals. The studied characteristics include Zipf's law of abbreviation (Zipf, 1949; Chaabouni et al., 2019; Ueda and Washio, 2021; Ueda and Taniguchi, 2024), word boundaries (Harris, 1955; Ueda et al., 2023; Ueda and Taniguchi, 2024), trade-off between word-order and case-marking (Comrie, 1989; Blake, 2001; Lian et al., 2023) and compositionality (Chaabouni et al., 2020; Rita et al., 2022b). On a more practical note, the language-like properties of induced protocols facilitate numerous applications. Mu et al. (2023) leverage emergent languages' superior functional expressivity for embodied control task. Yao et al. (2022) demonstrate the effectiveness of emergent languages in low-resource language modeling, and similar results are reported in machine translation (Li et al., 2020; Downey et al., 2023). Xu et al. (2022) show emergent languages' competitive as a representation learning method. Techniques for inducing compositionality in emergent languages (Zheng et al., 2024; Li and Bowling, 2019; Ren et al., 2020) find applications in improving generic neural networks' abilities (Ren et al., 2023; Zheng et al., 2024; Noukhovitch et al., 2023).

Environmental pressures for compositionality

Prerequisite conditions for emergence of compositionality are extensively studied. Kuciński et al. (2021) theoretically prove that compositional languages are more robust to message corruption and empirically verify that noisy channels facilitate compositionality. Several studies explore how capacity of communication channel (Lazaridou et al., 2018; Chaabouni et al., 2020) or capacity of neural agents (Resnick et al., 2020) affect compositionality. Cheng et al. (2023) observe that compositional languages are easier to imitate and suggest that imitability may also be a driving force for compositionality. Chaabouni et al. (2022) emphasize the task difficulty in terms of scale. Iterated learning (Smith et al., 2003; Li and Bowling, 2019; Ren et al., 2020) framework investigates the effects of language transmission across generations and finds that languages' learnability for newly created agents provides crucial pressure for compositional-

ity.

Community structures in emergent communication Our study on the one-to-many communication regime is closely related to a line of works that investigates the effects of community structures on emergent languages. [Harding Graesser et al. \(2019\)](#) explore how independently formed communities’ languages evolve when these communities start to interact with each other. [Kim and Oh \(2021\)](#) investigate the effects of different communication graphs on the languages’ properties. Several studies observe that naively increasing the population size does not yield more structured languages ([Chaabouni et al., 2022](#); [Kim and Oh, 2021](#)). [Rita et al. \(2022a\)](#) argue that different learning speeds in populations facilitate language structures. [Michel et al. \(2023\)](#) observe that limiting the communication graph with partitioning induces compositionality and generalization to unseen partners. However, all of these studies focus on one-to-one game play; hence, does not model the effects of one-to-many communication. [Chaabouni et al. \(2022\)](#) consider a simple voting mechanism of listeners only at inference time. [Li and Bowling \(2019\)](#) utilize naive message broadcasting when studying the effects of population size in iterated learning, but do not observe substantial improvements. [Yu et al. \(2022\)](#) employ message broadcasting in their work, but their main focus lies on adversarial aspects of communication, and they do not analyze compositionality of induced languages.

3 One-to-many communication game

We analyze emergent languages of agents playing a variant of Lewis reconstruction game ([Lewis, 1969](#)). The process of the game is as follows. Speaker π_θ observes an object $x \in \mathcal{X}^K$ and produces a message $m \sim \pi_\theta(\cdot | x)$ describing the object. An object contains K attributes and each attribute can take one of $|\mathcal{X}|$ possible values. A message $m \in \mathcal{W}^T$ is a sequence of symbols of fixed length T and each symbol belongs to vocabulary \mathcal{W} . The game contains a set of N listeners $\{\pi_{\phi_i}\}_{i=1}^N$. Each listener π_{ϕ_i} is concerned with K_i attributes, where $1 \leq K_i \leq K$. Let $x_i \in \mathcal{X}^{K_i}$ denote the K_i attributes’ values the listener π_{ϕ_i} is concerned with in object x , e.g., if $K_i = K$, then $x_i = x$.

For each round of game play, the set of listeners are randomly partitioned into M groups $\{\mathcal{G}_j\}_{j=1}^M$ such that $\cup_{j=1}^M \mathcal{G}_j = \{\pi_{\phi_i}\}_{i=1}^N$ and $\cap_{j=1}^M \mathcal{G}_j = \emptyset$.

Upon receiving message m , listener π_{ϕ_i} outputs its prediction for the object as $\hat{x}_i \sim \pi_{\phi_i}(\cdot | m)$. Let $\mathcal{G}^{(i)}$ denote the indices of listeners in the group that listener π_{ϕ_i} belongs to. Listener π_{ϕ_i} receives a reward of 1 if all of the listeners’ predictions in its group are correct, i.e., $R_{L_i}(x) = 1$ if $\forall j \in \mathcal{G}^{(i)}, \hat{x}_j = x_j$ and 0 otherwise. The speaker receives the average reward of all listeners as a reward, which is equal to the fraction of successful listeners: $R_S(x) = \frac{1}{N} \sum_{i=1}^N R_{L_i}(x)$. See Appendix A for graphical illustrations.

4 Experimental setup

Dataset We represent each attribute’s value with one-hot encoding. The number of attributes, K , is set to 4, and the number of values, $|\mathcal{X}|$, is set to 10. We set aside 10% of all attribute-value combinations as test set and use the rest as train set.

Speaker architecture One-hot encoded object x passes through a linear layer and initializes a single-layer GRU ([Cho et al., 2014](#)) of hidden size 500. It recurrently processes the input in total of $T = 5$ time steps. In each time step, the output is fed to a linear layer and then goes through Softmax activation to produce a vocabulary distribution of dimension $|\mathcal{W}| = 10$.

Listener architecture A listener π_{ϕ_i} is a single-layer GRU of hidden size 500. The listener sequentially processes the speaker’s message m , the last output of which is then passed to K_i linear layers corresponding to the number of attributes the listener is interested in. They each go through Softmax activation and produce a distribution of size $|\mathcal{X}|$ corresponding to the number of possible values an attribute can take.

Optimization We maximize each of the listeners’ and speaker’s expected reward with the REINFORCE algorithm ([Williams, 1992](#)). The expected reward for listener π_{ϕ_i} is written as $J_{L_i}(\phi_i) = \mathbb{E}_{x \sim p} \mathbb{E}_{m \sim \pi_\theta(\cdot | x)} R_{L_i}(x)$ and that of the speaker is written as $J_S(\theta) = \mathbb{E}_{x \sim p} \mathbb{E}_{m \sim \pi_\theta(\cdot | x)} R_S(x)$, where p denotes the uniform distribution over \mathcal{X}^K . We also utilize entropy regularization for the speaker to facilitate exploration and cross-entropy loss from listeners for stable training. We apply early stopping if the train set accuracy reaches 99%. Full description of the setup is in Appendix B.

Reporting We report average scores over 10 random seeds unless stated otherwise. Throughout

the paper, we use error bars to indicate 95% confidence interval and \pm to denote standard deviation. **Bold** and underline indicate the best and second best results.

5 Evaluation metrics

Topographic similarity (TopSim) Let $D_{\text{obj}} : \mathcal{X}^K \times \mathcal{X}^K \rightarrow \mathbb{R}^+$ and $D_{\text{msg}} : \mathcal{W}^T \times \mathcal{W}^T \rightarrow \mathbb{R}^+$ be distance measures over the objects and messages, respectively. Topographic similarity (Brighton and Kirby, 2006) is Spearman’s rank correlation of D_{obj} and D_{msg} over the joint uniform object, message distribution. High TopSim scores indicate that similar objects are mapped to similar messages. For D_{obj} and D_{msg} , we use cosine distance and Levenshtein distance (Levenshtein, 1965), respectively.

Positional disentanglement (PosDis) Let m_i denote the symbol in the i -th position of message m . Let a_1^i denote the attribute that has the highest mutual information with m_i , i.e., $a_1^i = \arg \max_a \mathcal{I}(m_i; a)$. Similarly, let a_2^i denote the attribute that has the second highest mutual information with m_i , i.e., $a_2^i = \arg \max_{a \neq a_1^i} \mathcal{I}(m_i; a)$. Positional disentanglement (Chaabouni et al., 2020) is equal to $\frac{1}{T} \sum_{i=1}^T \frac{\mathcal{I}(m_i; a_1^i) - \mathcal{I}(m_i; a_2^i)}{\mathcal{H}(m_i)}$, where $\mathcal{H}(m_i)$ denotes the entropy of the i -th position in messages. This measures the degree to which a single position is responsible for conveying information about an attribute.

Bag-of-symbols disentanglement (BosDis) Let n_i denote the number of occurrences of the i -th symbol from vocabulary \mathcal{W} in a message. Other notations follow from positional disentanglement. Bag-of-symbols disentanglement (Chaabouni et al., 2020) is equal to $\frac{1}{|\mathcal{W}|} \sum_{i=1}^{|\mathcal{W}|} \frac{\mathcal{I}(n_i; a_1^i) - \mathcal{I}(n_i; a_2^i)}{\mathcal{H}(n_i)}$. This measures how much a symbol univocally refers to an attribute. We provide a detailed description of PosDis and BosDis in Appendix E.

Compositional generalization Compositional generalization is the average task success rate on unseen attribute combinations. This is calculated using the test set without regard to the group.

6 Experiments

6.1 Does naive one-to-many communication enhance compositionality of languages?

Setup In naive one-to-many communication regime, all listeners share the same interests, and

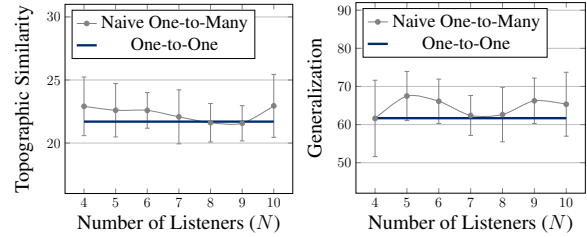


Figure 1: Language properties under varying numbers of listeners in naive one-to-many communication regime.

there is no coordination required among the listeners. More specifically, the number of attributes each listener is interested in is identical to the number of attributes the speaker observes ($K_i = K$), and each group contains only a single listener ($|\mathcal{G}_j| = 1$).

Naive message broadcasting does not improve compositionality Figure 1 compares languages from naive one-to-many communication regime with varying numbers of listeners (N) against the single-listener one-to-one communication regime ($N = 1$). While some of the cases exhibit improvements, none of the differences are statistically significant (two-tailed t -test with $p = 0.05$). The results suggest that simply broadcasting a message does not introduce a meaningful pressure on language emergence.

6.2 How do listeners of different interests affect language properties?

Setup We devise three kinds of listener formations for this experiment. The *partial-interest* formation contains $\binom{K}{K_i}$ listeners that are only concerned with K_i attributes. Each of $\binom{K}{K_i}$ listeners’ interests are distinct attribute combinations. The *mixed-interest* formation is the same as the partial-interest formation except that it contains one additional listener that is concerned with all of the K attributes. The *full-interest* formation contains $1 + \binom{K}{K_i}$ listeners all of which are interested in all of the K attributes. As there is no coordination required, each group \mathcal{G}_i contains a single listener. The test set accuracy is calculated only with the listeners that are interested in all of the attributes. Appendix G contains experiments on a larger setup.

Readability pressure from different interests facilitates more structured languages In Figure 2a, we observe a trend that the more the listeners can disregard other parts of a message that they are

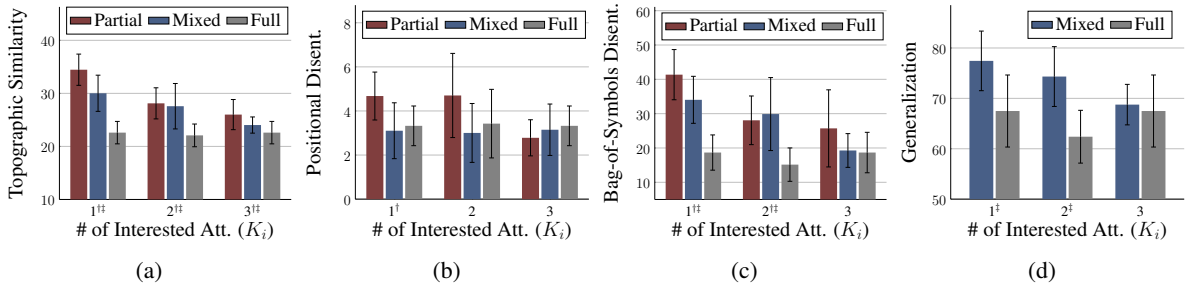


Figure 2: Comparison of language properties in listeners of different interests regime. † and ‡ denote statistically significant differences for the partial- and mixed-interest formations compared to the full-interest formation (one-tailed t -test with $p < 0.05$).

not concerned with, the more compositional the languages tend to be. The formations with smaller number of interested attributes (K_i) exhibit higher TopSim, and the partial-interest formation’s languages tend to exhibit higher TopSim compared to the mixed-interest formation. Languages from the two formations are more compositional than those of a similarly sized full-interest formation. In Figure 2d, we observe a similar trend for compositional generalization ability. We hypothesize that listeners of different interests prefer more structured languages, so that they can more easily infer the attributes of interest from a message without needing to understand other details that are not related to their interests. We confirm that the results do not stem from the relative easiness of the task in Appendix F.

Listeners of different interests prefer symbolic structures rather than position We analyze what kinds of language structures are promoted by listeners of different interests. One possible structure is to denote each attribute within a certain position of a message. However, we do not observe such positional structures in regard to the number of interested attributes from Figure 2b. Another possible strategy is to associate the number of occurrences of a certain symbol to an attribute. In Figure 2c, we observe a clear trend that listeners of different interests prefer this kind of association when listeners are interested in smaller subsets of the whole attributes.

Languages from listeners of different interests are easier to learn We test if listeners of different interests indeed facilitate more structured, hence easier to learn languages. We take languages from the partial-interest formation with the number of interested attributes set to one ($K_i = 1$) and the full-interest formation of equal size. We randomly

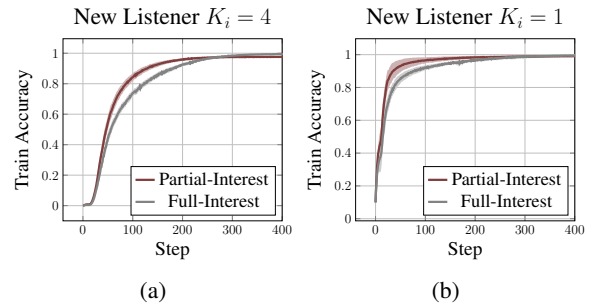


Figure 3: Learnability comparison in different interests regime. Shades indicate one standard deviation.

initialize new listeners of two different interests; one is only interested in one randomly sampled attribute ($K_i = 1$), and the other is interested in all of the four attributes ($K_i = 4$). We train these listeners by letting them play the game with the frozen speakers of respective languages. In Figure 3, we observe that in both cases the languages from the partial-interest formation are easier to learn.

6.3 How does coordination pressure affect language properties?

Setup We construct 50 listeners of the same interests ($K_i = K$). For each round of game play, the listeners are randomly split into equally sized groups. We explore the effects of coordination pressure in terms of group size ($|\mathcal{G}_j|$). A larger group size forces more listeners to be simultaneously successful at understanding the speaker’s message. The test accuracy is calculated by taking average of all listeners’ success rates regardless of the group.

Coordination pressure amplifies preference for compositionality In Figure 4a, we observe increase in TopSim as sufficient degrees of coordination pressure are introduced to the game. Increases corresponding to group sizes greater than 2 are statistically significant. TopSim increases with the

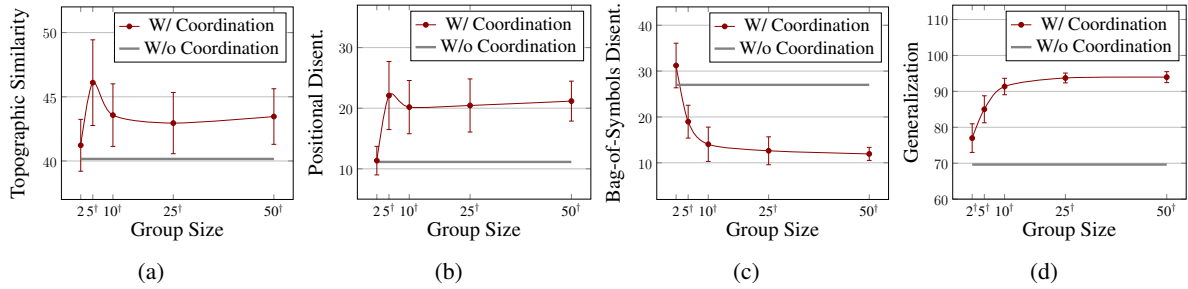


Figure 4: Comparison of language properties under varying degrees of coordination pressure. [†] denotes statistically significant difference compared to group size of 1 (one-tailed t -test with $p < 0.05$, averaged over 20 random seeds).

group size up to 5, then tends to decrease at larger group sizes. A more steady increase is observed in test set accuracy in Figure 4d. We hypothesize that the coordination pressure amplifies the degree of preference for the language’s compositionality from the listeners, as it requires the listeners to have a simultaneously shared understanding of a message.

Coordination pressure induces position-wise structures rather than symbols In Figure 4b, we observe increase in PosDis when sufficient degrees of coordination pressure are injected to the game, suggesting that coordination pressures induce more position-wise structures. A reverse trend is observed in BosDis in Figure 4c. The emergent languages under coordinate pressure tend to rely less on the number of occurrences of a symbol when determining an attribute’s value. The results indicate that to effectively express more complicated concepts (larger number of attributes) position-wise structures are preferred.

6.4 Coordination pressure in relation to iterated learning framework

Iterated learning Iterated learning framework (Smith et al., 2003) simulates languages’ transmission across generations. Li and Bowling (2019) find that periodically resetting listener’s parameters forces the speaker to develop languages that are easier to teach and more compositional. In their experiments with populations of listeners, the authors hypothesize that resetting one listener at a time in a staggered manner instead of resetting them all at once could yield more structured languages as the population would contain more diverse listeners with varying degrees of experience. However, they observe that simultaneously resetting all of the listeners at the same time yield more compositional languages compared to the staggered reset regime

and suggest that improvements of iterated learning can be attributed to the abrupt changes induced from simultaneous resets of listeners.

Setup We conduct a small-scale experiment with two listeners to explore how coordination pressure impacts languages in iterated learning. We consider three different listener reset regimes. In simultaneous reset regime, we reset all listeners every 1,000 epochs. Staggered reset regime resets one listener at epochs $\{500, 1500, 2500, \dots\}$ and the other listener at epochs $\{1000, 2000, 3000, \dots\}$. No-reset regime does not perform any listener resets. We also consider a single-listener system under no-reset and simultaneous reset regimes. We train the agents for 10,000 epochs.

Coordination pressure enhances compositionality in more realistic population dynamics In Table 1, we compare the single-listener system to the two-listener systems with and without coordination pressure (group size of 2 and 1, respectively). Without coordination pressure, staggered reset produces less compositional languages compared to the simultaneous reset regime. When coordination pressure is introduced to the game, staggered reset exhibits comparable compositionality to the simultaneous reset regime. Under coordination pressure, a newly reset listener may affect the game performance of another listener as they are required to coordinate. We hypothesize that this introduces abruptness to the system and that the language properties of the staggered reset regime are enhanced by this added abruptness. The experimental results suggest that coordination pressure is an important dimension in iterated learning framework as it can enhance the language structure in a more realistic setup of population dynamics. It is often less likely that the entire population is replaced at once than that it is to undergo a gradual change, as is the case in the staggered reset regime.

Metric	Single Listener ($N = 1$)		Without Coordination ($N = 2$)			With Coordination ($N = 2$)		
	No-Reset	Simultaneous	No-Reset	Simultaneous	Staggered	No-Reset	Simultaneous	Staggered
TopSim	25.97±3.0	33.36±3.6	25.77±3.0	34.62±3.4	29.52±3.1	27.60±4.0	34.79±6.2	34.25±5.0 [†]
Generalization	61.82±8.5	84.44±6.2	61.13±9.4	81.74±6.6	69.77±9.0	69.08±15.3	85.08±10.6	84.25±8.6 [†]

Table 1: Effects of coordination pressure on emergent languages in iterated learning environment. [†] indicates statistically significant difference from the corresponding reset regime in the middle block (one-tailed t -test with $p < 0.05$).

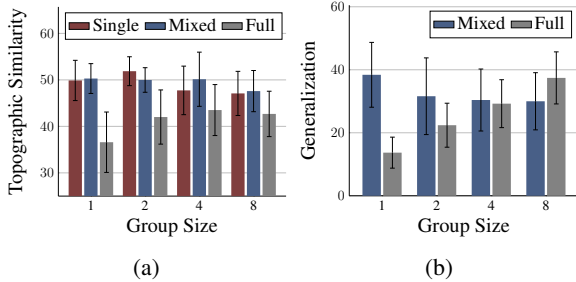


Figure 5: Comparison of language properties in general one-to-many communication regime.

6.5 Listeners of different interests under coordination pressure

We explore how the readability pressure from listeners of different interests and coordination pressure interact with each other in language emergence.

Setup To observe the interaction in a more granular scale, we increase the number of attributes, K , to 8. The number of possible values, $|\mathcal{X}|$, is in turn reduced to 2 to stabilize training. The length of messages, T , is also reduced to 3. We construct three kinds of listener formations. The *single-interest* formation contains eight listeners that are interested in each of the eight attributes ($K_i = 1$). The *mixed-interest* formation is the same as the single-interest formation but contains one additional listener that is interested in all of the attributes. The *full-interest* formation contains nine listeners that are interested in all of the eight attributes ($K_i = K$). We test how these listener formations behave under varying degrees of coordination pressure expressed by group sizes of 1, 2, 4, 8. We note that the mixed-interest formation leaves one single-listener group at group sizes greater than 1 as it contains one additional listener.

Readability and coordination pressures can clash in different directions In Figure 5a, a slight decreasing tendency in TopSim is observed in the single- and mixed-interest formations as group size is increased. This is in contrast to the case of the full-interest formation where TopSim val-

ues exhibit an increasing tendency with increase in the group size. We observe a similar trend in Figure 5b for test set accuracy. These observations suggest that coordination pressure can counteract readability pressure induced from listeners of different interests regime. Results in §6.2 show that the more the listeners can disregard other parts of a message that they are not concerned with, which is manifested in lower numbers of interested attributes, the more compositional the languages tend to be. We hypothesize that coordination pressure may conflict with readability pressure as coordination among listeners forces larger numbers of attributes to be predicted correctly at the same time. From the speaker’s point of view, these listeners act as though they are interested in a larger number of attributes.

7 Experiments with raw images

We expand our study to more realistic scenarios employing datasets that consist of raw pixel images.

7.1 Listeners of different interests with raw pixel data

Experimental setup We explore the effects of readability pressure introduced by listeners of different interests in a more realistic setup with 3dshapes dataset (Kim and Mnih, 2018). The dataset contains images of 3D shapes. Each image is characterized by 6 attributes such as the object’s color and shape. We sample 4 values from each of these 6 attributes and perform the same experiment as in §6.2. Full description of the experimental setup is in Appendix C.

Results Overall, we observe similar trends to those of the attribute-value dataset, suggesting that the findings in §6.2 hold in more complex environments. In Figure 6a, we find that smaller numbers of attributes of interest yield more compositional languages, and Figure 6d shows that they exhibit stronger generalization ability. We obtain more

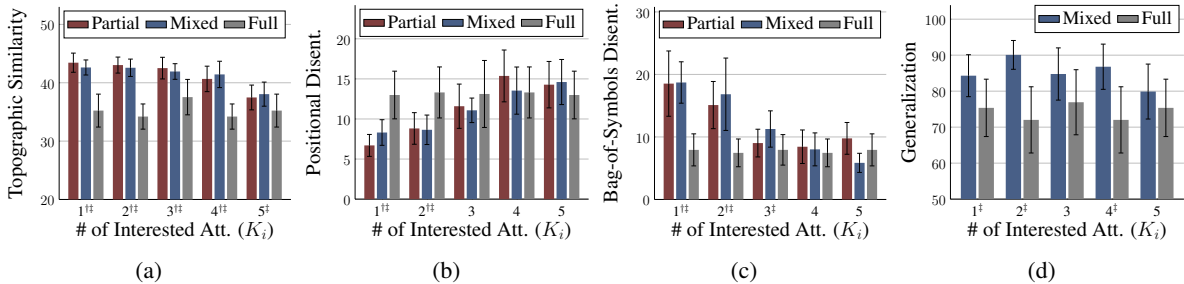


Figure 6: Comparison of language properties in listeners of different interests regime on 3dshapes dataset. \dagger and \ddagger denote statistically significant differences for the partial- and mixed-interest formations compared to the full-interest formation (one-tailed t -test with $p < 0.05$).

pronounced effects in terms of symbol- or position-wise structures of emergent languages. There is a clear tendency that smaller number of interested attributes produce languages that are less reliant on positional structures of messages as can be seen in Figure 6b. In Figure 6c, we also observe the tendency to denote an attribute with number of occurrences of a symbol in listeners of different interests regime. We evaluate the degree of association between values of attributes and symbols in Appendix H.

7.2 Coordination pressure in large scale image discrimination game

Discrimination game We explore the effects of coordination pressure in a large-scale image discrimination game with ImageNet dataset (Rusakovsky et al., 2015). The rules of the game are as follows. The speaker observes the target image x and sends a message m containing descriptions of the image to a set of listeners $\{\pi_{\phi_i}\}_{i=1}^N$. A listener π_{ϕ_i} is tasked to determine which one is the target among its context \mathcal{C}_i containing other images and rewarded if all of the listeners in the group it belongs to correctly predict the target.

Scramble resistance (ScrRes) Let m' denote a randomly permuted version of message m and $\pi_{\phi_i}(x | m, \mathcal{C}_i)$ denote the probability assigned to the target object x by listener π_{ϕ_i} given message m and context \mathcal{C}_i . Scramble resistance (Bernard and Mickus, 2023) is calculated as $\frac{\min(\pi_{\phi_i}(x|m, \mathcal{C}_i), \pi_{\phi_i}(x|m', \mathcal{C}_i))}{\pi_{\phi_i}(x|m, \mathcal{C}_i)}$. A high scramble resistance score indicates that the language is less affected by positional perturbations.

Experimental setup We use representations of images processed by a ResNet-50 (He et al., 2016) encoder pretrained on ImageNet with BYOL (Grill et al., 2020) as in Chaabouni et al. (2022); Michel

Group Size	Task Success Rate		Compositionality	
	Test (OOD)	Val (ID)	TopSim	ScrRes
1	86.92±0.5	93.03±3.3	20.59±1.4	35.23±2.3
2	88.06±1.0 \dagger	93.66±4.1	19.40±1.6	35.75±4.7
5	87.76±0.7 \dagger	93.46±3.7	20.38±1.8	35.46±3.7
10	89.69±0.3 \ddagger	94.55±1.7 \dagger	21.68±2.3	30.55±1.6 \dagger

Table 2: Results on image discrimination game. \dagger denotes statistically significant difference compared to group size of 1 (one-tailed t -test with $p < 0.05$).

et al. (2023). The context size $|\mathcal{C}_i|$ is set to 100 for all listeners. We use the train, validation, test splits from Chaabouni et al. (2022). We set aside 10% of the classes in the dataset as in-distribution (ID) classes and the rest as out-of-distribution (OOD) classes. We perform training and validation with ID samples in each split and evaluation with the test set containing only OOD samples. TopSim is calculated with respect to the image’s representations using cosine distance. We construct 10 listeners and observe the effects of coordination pressure under varying group sizes. Full description of the experimental setup is in Appendix D.

Results We report the accuracies on each split as well as compositionality measures in Table 2. We observe that coordination pressure induces stronger generalization ability in both OOD and ID samples. Group size of 10 exhibits the highest generalization ability. We do not observe a clear correlation in TopSim and generalization ability. Prior works (Chaabouni et al., 2022; Michel et al., 2023) also report that TopSim does not correlate with generalization ability and suggest that it may be an inadequate measure of compositionality for complex data forms. A lower value of ScrRes at group size of 10 indicates that coordination pressure also induces more position-wise structures in languages in more complex setups, but its effect at lower group

sizes is not as clear.

8 Conclusion

This work investigates how one-to-many communication affects language emergence. We find that one-to-many communication introduces two aspects of communication that facilitate emergence of compositionality. First, listeners of different interests exert readability pressure. This forces the language to be more structured as listeners prefer messages that do not require understanding of other aspects unrelated to the attributes of interest. Second, coordination among listeners forces languages to take more structured forms as it makes languages easier to be simultaneously understood by multiple listeners. Additionally, we find that coordination promotes emergence of compositionality in more realistic population dynamics. We verify that our findings hold in more complex environments with experiments on raw image data. Our work sheds light on the importance of one-to-many communication in the emergent communication field.

Limitations

Task complexity This work analyzes emergent languages with basic attribute-values and image datasets. While these datasets are widely employed in the emergent communication community and permit a detailed analysis of compositionality, they lack the complexities of real-world environments. Recent studies propose various tasks that require more abstract reasoning (Guo et al., 2023; Zhou et al., 2024; Mihai and Hare, 2021; Patel et al., 2021). Future work may explore how our findings apply in more complex task scenarios.

Compositionality measures Measuring compositionality of a language is a challenging task, and existing metrics are known to measure only crude forms of compositionality (Korbak et al., 2020). This hinders scaling up task complexity as existing measures may not be adequate for more complex forms of compositionality (Chaabouni et al., 2022; Michel et al., 2023). This also limits assessing the implications and impacts of our work as existing metrics may fail to reflect more nuanced and complex aspects of compositionality of human languages.

Complex communication structures This study sets a basic one-to-many communication of a single speaker and the speaker’s message is broadcast

to all listeners in the system. However, more complex communication structures are possible. There could be multiple speakers and a speaker’s message may be relayed to only certain portions of the listeners. The effects of population size (Rita et al., 2022a; Michel et al., 2023) and more complex communication graphs (Kim and Oh, 2021; Harding Graesser et al., 2019; Michel et al., 2023) could be further explored. On the coordination side, instead of forming new groups for each game play, longer listener group formation frequency could be explored. We also note that the effects of skewed interests of listeners are not explored in this work as we simply utilized all combinations of interests.

Exploration of applications Our work does not explore immediate application areas of the findings. However, the emergent communication field has demonstrated numerous application possibilities in diverse domains. Some of these find applications in improving foundation models (Noukhovitch et al., 2023; Zheng et al., 2024). It may be an interesting research direction to investigate our findings in relation to alignment of large language models (Ouyang et al., 2022; Rafailov et al., 2023) as human preferences can be decomposed into multiple attributes (Lou et al., 2024), e.g., helpfulness, politeness, etc. Our findings suggest that devising separate preference models, each concerning a certain preference aspect, could be beneficial for compositional generalization in terms of these preferences. As for the coordination pressure, multiple preference models of different value systems could be explored for simultaneously satisfying a wide range of users of varying cultural backgrounds.

Theoretical analysis Through extensive experiments, we empirically verify that listeners of different interests and coordination among listeners play crucial roles in emergence of compositionality. However, more fine-grained analyses of the process would enhance the understanding of these factors and facilitate application possibilities. One could theoretically analyze the processing efforts required for listeners of different interests are indeed lower when the language is more compositional, or theoretically validate that the chances of any two listeners to stumble upon the same protocol are higher when the language is compositional.

We provide further discussions on limitations in Appendix I.

Acknowledgements

We thank the anonymous reviewers for their insightful comments and constructive feedback, which have significantly improved this work.

Figures 7, 8 contain icons provided by Flaticon.

References

- Timothée Bernard and Timothee Mickus. 2023. [So many design choices: Improving and interpreting neural agent communication in signaling games](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 8399–8413, Toronto, Canada. Association for Computational Linguistics.
- B.J. Blake. 2001. *Case*. Cambridge Textbooks in Linguistics. Cambridge University Press.
- Diane Bouchacourt and Marco Baroni. 2018. [How agents see things: On visual representations in an emergent language game](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 981–985, Brussels, Belgium. Association for Computational Linguistics.
- Henry Brighton and Simon Kirby. 2006. [Understanding Linguistic Evolution by Visualizing the Emergence of Topographic Mappings](#). *Artificial Life*, 12(2):229–242.
- T. Briscoe. 2002. *Linguistic Evolution through Language Acquisition*. Cambridge University Press.
- Boaz Carmeli, Yonatan Belinkov, and Ron Meir. 2024. [Concept-best-matching: Evaluating compositionality in emergent communication](#). In *Findings of the Association for Computational Linguistics ACL 2024*, pages 3186–3194, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.
- Rahma Chaabouni, Eugene Kharitonov, Diane Bouchacourt, Emmanuel Dupoux, and Marco Baroni. 2020. [Compositionality and generalization in emergent languages](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4427–4442, Online. Association for Computational Linguistics.
- Rahma Chaabouni, Eugene Kharitonov, Emmanuel Dupoux, and Marco Baroni. 2019. *Anti-efficient encoding in emergent communication*. Curran Associates Inc., Red Hook, NY, USA.
- Rahma Chaabouni, Florian Strub, Florent Althé, Eugene Tarassov, Corentin Tallec, Elnaz Davoodi, Kory Wallace Mathewson, Olivier Tieleman, Angeliki Lazaridou, and Bilal Piot. 2022. [Emergent communication at scale](#). In *International Conference on Learning Representations*.
- Emily Cheng, Mathieu Rita, and Thierry Poibeau. 2023. [On the correspondence between compositionality and imitation in emergent neural communication](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 12432–12447, Toronto, Canada. Association for Computational Linguistics.
- Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. [Learning phrase representations using RNN encoder–decoder for statistical machine translation](#). In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1724–1734, Doha, Qatar. Association for Computational Linguistics.
- Edward Choi, Angeliki Lazaridou, and Nando de Freitas. 2018. [Multi-agent compositional communication learning from raw visual input](#). In *International Conference on Learning Representations*.
- B. Comrie. 1989. *Language Universals and Linguistic Typology: Syntax and Morphology*. University of Chicago Press.
- Donald Davidson. 1965. Theories of meaning and learnable languages. In Yehoshua Bar-Hillel, editor, *Proceedings of the 1964 International Congress for Logic, Methodology, and Philosophy of Science*, pages 383–394. North-Holland.
- C.m. Downey, Xuhui Zhou, Zeyu Liu, and Shane Steinert-Threlkeld. 2023. [Learning to translate by learning to communicate](#). In *Proceedings of the 3rd Workshop on Multi-lingual Representation Learning (MRL)*, pages 218–238, Singapore. Association for Computational Linguistics.
- J.H. Greenberg. 1961. *Universals of Language*. MIT.
- Jean-Bastien Grill, Florian Strub, Florent Althé, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, Bilal Piot, koray kavukcuoglu, Remi Munos, and Michal Valko. 2020. [Bootstrap your own latent - a new approach to self-supervised learning](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 21271–21284. Curran Associates, Inc.
- Yuxuan Guo, Yifan Hao, Rui Zhang, Enshuai Zhou, Zidong Du, Xishan Zhang, Xinkai Song, Yuanbo Wen, Yongwei Zhao, Xuehai Zhou, Jiaming Guo, Qi Yi, Shaohui Peng, Di Huang, Ruizhi Chen, Qi Guo, and Yunji Chen. 2023. [Emergent communication for rules reasoning](#). In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Laura Harding Graesser, Kyunghyun Cho, and Douwe Kiela. 2019. [Emergent linguistic phenomena in multi-agent communication games](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3700–3710, Hong Kong, China. Association for Computational Linguistics.

- Zellig S. Harris. 1955. [From phoneme to morpheme](#). *Language*, 31(2):190–222.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. [Deep residual learning for image recognition](#). In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778.
- Theo M.V. Janssen and Barbara H. Partee. 1997. [Chapter 7 - compositionality](#). In Johan van Benthem and Alice ter Meulen, editors, *Handbook of Logic and Language*, pages 417–473. North-Holland, Amsterdam.
- Eugene Kharitonov, Roberto Dessì, Rahma Chaabouni, Diane Bouchacourt, and Marco Baroni. 2021. EGG: a toolkit for research on Emergence of lanGuage in Games. <https://github.com/facebookresearch/EGG>.
- Hyunjik Kim and Andriy Mnih. 2018. [Disentangling by factorising](#). In *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 2649–2658. PMLR.
- Jooyeon Kim and Alice Oh. 2021. [Emergent communication under varying sizes and connectivities](#). In *Advances in Neural Information Processing Systems*.
- Diederik P. Kingma and Jimmy Ba. 2015. [Adam: A method for stochastic optimization](#). In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.
- Tomasz Korbak, Julian Zubek, and Joanna Rączaszek-Leonardi. 2020. [Measuring non-trivial compositionality in emergent communication](#). In *4th Workshop on Emergent Communication (Talking to Strangers: Zero-Shot Emergent Communication) at NeurIPS 2020*.
- Łukasz Kuciński, Tomasz Korbak, Paweł Koł odziej, and Piotr Miłoś. 2021. [Catalytic role of noise and necessity of inductive biases in the emergence of compositional communication](#). In *Advances in Neural Information Processing Systems*, volume 34, pages 23075–23088. Curran Associates, Inc.
- Angeliki Lazaridou and Marco Baroni. 2020. [Emergent multi-agent communication in the deep learning era](#). *Preprint*, arXiv:2006.02419.
- Angeliki Lazaridou, Karl Moritz Hermann, Karl Tuyls, and Stephen Clark. 2018. [Emergence of linguistic communication from referential games with symbolic and pixel input](#). In *International Conference on Learning Representations*.
- Angeliki Lazaridou, Alexander Peysakhovich, and Marco Baroni. 2017. [Multi-agent cooperation and the emergence of \(natural\) language](#). In *International Conference on Learning Representations*.
- Vladimir I. Levenshtein. 1965. Binary codes capable of correcting deletions, insertions, and reversals. *Soviet physics. Doklady*, 10:707–710.
- David Kellogg Lewis. 1969. *Convention: A Philosophical Study*. Wiley-Blackwell, Cambridge, MA, USA.
- Fushan Li and Michael Bowling. 2019. [Ease-of-teaching and language structure from emergent communication](#). In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.
- Yaoyiran Li, Edoardo Maria Ponti, Ivan Vulić, and Anna Korhonen. 2020. Emergent communication pretraining for few-shot machine translation. In *Proceedings of the 28th International Conference on Computational Linguistics*.
- Yuchen Lian, Arianna Bisazza, and Tessa Verhoef. 2023. [Communication drives the emergence of language universals in neural agents: Evidence from the word-order/case-marking trade-off](#). *Transactions of the Association for Computational Linguistics*, 11:1033–1047.
- Xingzhou Lou, Junge Zhang, Jian Xie, Lifeng Liu, Dong Yan, and Kaiqi Huang. 2024. [Spo: Multi-dimensional preference sequential alignment with implicit reward modeling](#). *Preprint*, arXiv:2405.12739.
- Paul Michel, Mathieu Rita, Kory Wallace Mathewson, Olivier Tieleman, and Angeliki Lazaridou. 2023. [Revisiting populations in multi-agent communication](#). In *The Eleventh International Conference on Learning Representations*.
- Daniela Mihai and Jonathon Hare. 2021. [Learning to draw: Emergent communication through sketching](#).
- Jesse Mu and Noah Goodman. 2021. [Emergent communication of generalizations](#). In *Advances in Neural Information Processing Systems*.
- Yao(Mark) Mu, Shunyu Yao, Mingyu Ding, Ping Luo, and Chuang Gan. 2023. [Ec2: Emergent communication for embodied control](#). *The IEEE / CVF Computer Vision and Pattern Recognition Conference (CVPR)*.
- Michael Noukhovitch, Samuel Lavoie, Florian Strub, and Aaron C Courville. 2023. [Language model alignment with elastic reset](#). In *Advances in Neural Information Processing Systems*, volume 36, pages 3439–3461. Curran Associates, Inc.
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. [Training language models to follow instructions with human feedback](#). *Preprint*, arXiv:2203.02155.

- Shivansh Patel, Saim Wani, Unnat Jain, Alexander Schwing, Svetlana Lazebnik, Manolis Savva, and Angel X. Chang. 2021. Interpretation of emergent communication in heterogeneous collaborative embodied agents. In *ICCV*.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. [Direct preference optimization: Your language model is secretly a reward model](#). In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Limor Raviv, Antje Meyer, and Shiri Lev-Ari. 2019. [Larger communities create more systematic languages](#). *Proceedings of the Royal Society B: Biological Sciences*, 286:20191262.
- Yi Ren, Shangmin Guo, Matthieu Labeau, Shay B. Cohen, and Simon Kirby. 2020. [Compositional languages emerge in a neural iterated learning model](#). In *International Conference on Learning Representations*.
- Yi Ren, Samuel Lavoie, Mikhail Galkin, Danica J. Sutherland, and Aaron Courville. 2023. [Improving compositional generalization using iterated learning and simplicial embeddings](#). In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Cinjon Resnick, Abhinav Gupta, Jakob Foerster, Andrew M. Dai, and Kyunghyun Cho. 2020. [Capacity, bandwidth, and compositionality in emergent language learning](#). In *Proceedings of the 19th International Conference on Autonomous Agents and Multi-Agent Systems, AAMAS '20*, page 1125–1133, Richland, SC. International Foundation for Autonomous Agents and Multiagent Systems.
- Mathieu Rita, Florian Strub, Jean-Bastien Grill, Olivier Pietquin, and Emmanuel Dupoux. 2022a. [On the role of population heterogeneity in emergent communication](#). In *International Conference on Learning Representations*.
- Mathieu Rita, Corentin Tallec, Paul Michel, Jean-Bastien Grill, Olivier Pietquin, Emmanuel Dupoux, and Florian Strub. 2022b. [Emergent communication: Generalization and overfitting in lewis games](#). In *Advances in Neural Information Processing Systems*.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. 2015. [ImageNet Large Scale Visual Recognition Challenge](#). *International Journal of Computer Vision (IJCV)*, 115(3):211–252.
- Jürgen Schmidhuber. 2015. [Deep learning in neural networks: An overview](#). *Neural Networks*, 61:85–117.
- Kenny Smith, Simon Kirby, and Henry Brighton. 2003. [Iterated learning: A framework for the emergence of language](#). *Artificial life*, 9:371–86.
- Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. 1999. [Policy gradient methods for reinforcement learning with function approximation](#). In *Advances in Neural Information Processing Systems*, volume 12. MIT Press.
- Olivier Tieleman, Angeliki Lazaridou, Shibl Mourad, Charles Blundell, and Doina Precup. 2019. Shaping representations through communication: community size effect in artificial learning systems. *Visually Grounded Interaction and Language (ViGIL) Workshop*.
- Ryo Ueda, Taiga Ishii, and Yusuke Miyao. 2023. [On the word boundaries of emergent languages based on harris’s articulation scheme](#). In *The Eleventh International Conference on Learning Representations*.
- Ryo Ueda and Tadahiro Taniguchi. 2024. [Lewis’s signaling game as beta-VAE for natural word lengths and segments](#). In *The Twelfth International Conference on Learning Representations*.
- Ryo Ueda and Koki Washio. 2021. [On the relationship between Zipf’s law of abbreviation and interfering noise in emergent languages](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing: Student Research Workshop*, pages 60–70, Online. Association for Computational Linguistics.
- Ronald J Williams. 1992. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8:229–256.
- Zhenlin Xu, Marc Niethammer, and Colin Raffel. 2022. [Compositional generalization in unsupervised compositional representation learning: A study on disentanglement and emergent language](#). In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*.
- Shunyu Yao, Mo Yu, Yang Zhang, Karthik Narasimhan, Joshua Tenenbaum, and Chuang Gan. 2022. Linking emergent and natural languages via corpus transfer. In *International Conference on Learning Representations (ICLR)*.
- Dhara Yu, Jesse Mu, and Noah Goodman. 2022. [Emergent covert signaling in adversarial reference games](#). In *Emergent Communication Workshop at ICLR 2022*.
- Chenhao Zheng, Jieyu Zhang, Aniruddha Kembhavi, and Ranjay Krishna. 2024. Iterated learning improves compositionality in large vision-language models. *The IEEE / CVF Computer Vision and Pattern Recognition Conference (CVPR)*.
- Enshuai Zhou, Yifan Hao, Rui Zhang, Yuxuan Guo, Zidong Du, Xishan Zhang, Xinkai Song, Chao Wang, Xuehai Zhou, Jiaming Guo, Qi Yi, Shaohui Peng, Di Huang, Ruizhi Chen, Qi Guo, and Yunji

Chen. 2024. [Emergent communication for numerical concepts generalization](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(16):17609–17617.

George K. Zipf. 1949. *Human Behaviour and the Principle of Least Effort*. Addison-Wesley.

A Graphical illustration of one-to-many communication game

Figure 7 illustrates listeners of different interests in one-to-many communication game. The speaker’s message is broadcast to three listeners. These listeners each have their own distinct interests. The first listener is only interested in the color of the object, while the second listener is only interested in the shape of the object. The third listener is interested in both the color and the shape of the object. The predictions of these listeners reflect their interests, hence exclusively pertain to the attributes of interest.

Figure 8 illustrates coordination among four listeners. Each of the four listeners are assigned to a group of size 2. The speaker’s message is broadcast to the listeners, and each listener predicts the object’s attributes. Both listeners in the first group correctly predict the object’s attributes and the group is considered to be successful at the task. One of the listeners in the second group produces an incorrect prediction and this results in a failure of the task for the entire group.

B Experimental details

We utilize EGG framework (Kharitonov et al., 2021), which is available under MIT license. Speaker’s symbol embedding size is 5 and listeners’ symbol embedding size is 30. We use Adam optimizer (Kingma and Ba, 2015) with learning rate of 0.001. The batch size is set to 5120. We utilize REINFORCE with baseline (Sutton et al., 1999) where the baseline function is the average of the past rewards for the corresponding speaker or listener agent. We report compositionality metrics from the full dataset as they exhibit high correlation across splits, except for §6.5, where results on test set are reported as weak correlation is observed. We exclude a few runs that do not reach train accuracy of 99% in experiments not involving coordination pressure (§6.1, §6.2, §7.1). For experiments on coordination pressure with attribute-value datasets, training is performed for 5,000 epochs in §6.3 and 30,000 epochs in §6.5. At inference time,

messages are constructed by selecting the symbol that has been assigned the highest probability by the speaker at each time step. Experiments on raw pixel datasets follow the same setup unless otherwise specified.

Cross-entropy loss The training objective contains cross-entropy loss from listener to stabilize training process. The cross-entropy loss for listener π_{ϕ_i} is written as $-\frac{1}{K_i} \sum_{k=1}^{K_i} \log \pi_{\phi_i}(x_i^{(k)} | m)$, where $x_i^{(k)}$ refers to the k -th attribute in the object of interest x_i for the listener. For the speaker, listeners’ average cross-entropy loss is added to the reward after taking its negative. For the listeners, each listener’s own cross-entropy loss is added to the reward in a similar manner. In addition to that, we directly backpropagate the cross-entropy loss for each listener. Each cross-entropy loss term is multiplied by a scaling hyperparameter λ . This loss is coordination-agnostic in that it is not affected by the group members’ success or failure. In experiments that do not involve coordination pressure (§6.1, §6.2, §7.1), the value is set to 1.0. For experiments that involve coordination (§6.3, §6.5, §7.2), a lower value of 0.01 is utilized to better observe the effects of coordination pressure, except for iterated learning experiments in §6.4, where the value is set to 0.03 to offset the narrow training window induced from the parameter resets.

Entropy regularization We add entropy regularization term in the speaker’s symbol distribution to promote exploration. The magnitude of the regularization is controlled by a scaling hyperparameter γ , which is multiplied to the entropy term. γ is set to induce successful language emergence on the train set of each dataset. For the experiments in §6.1, §6.2, the value is set to 0.5. In the experiments with 3dshapes dataset, the value is set to 1.0. In the experiments that involve coordination (§6.3, §6.5, §7.2), the value is set to 0.01. In experiments on iterated learning, the value is set to 0.02.

C Experimental details on 3dshapes

We set the vocabulary size $|\mathcal{W}|$ to 6 and the length of messages T to 6. The batch size is set to 5,120. We stop training when the train accuracy reaches 99%. We run each experiment with 20 random seeds and report the average. The dataset is available under Apache-2.0 license.

Dataset construction An image is characterized by 6 attributes: object’s shape, object’s color, ob-

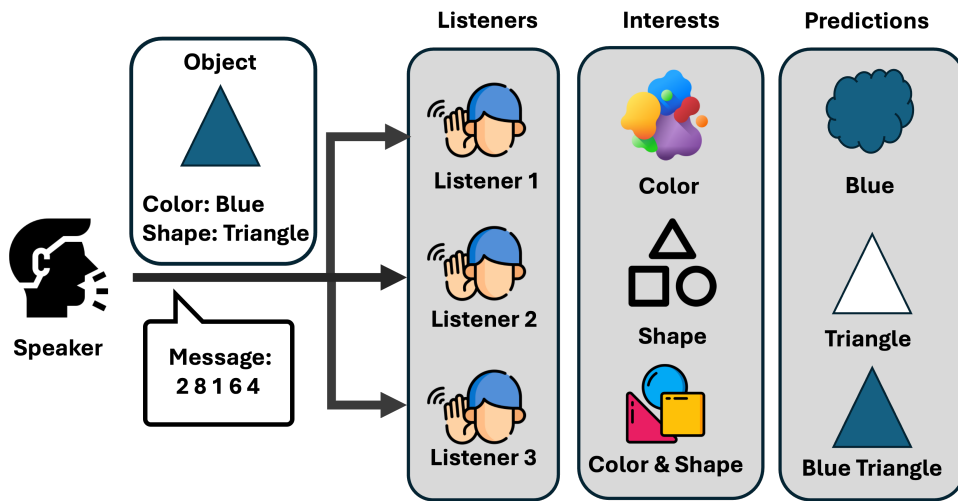


Figure 7: Illustration of listeners of different interests in one-to-many communication game. Each listener is interested in a different set of attributes and its predictions only pertain to the attributes of interest.

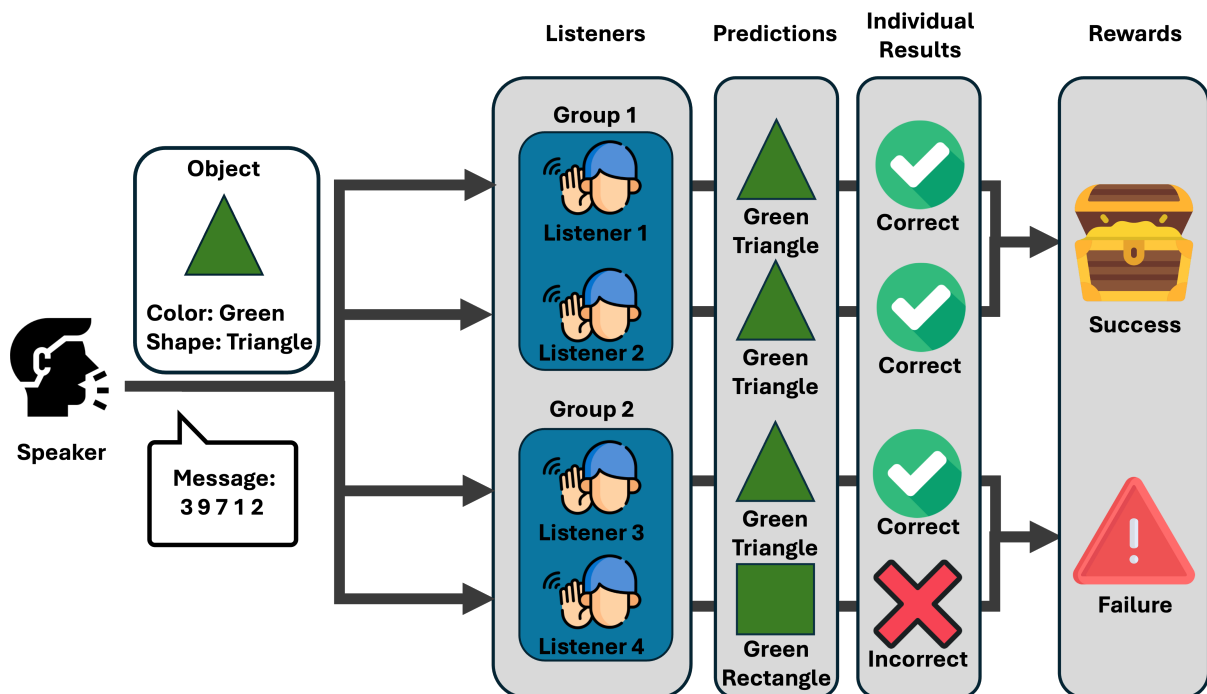


Figure 8: Illustration of coordination among listeners in one-to-many communication game. Listeners are split into groups and each listener is rewarded if and only if all of the listeners in the same group correctly predict the attributes.

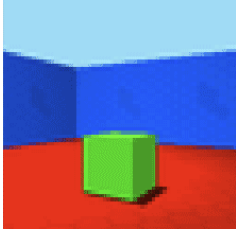


Figure 9: A sample of 3dshapes dataset.

ject’s size, color of the wall, color of the floor, and viewing orientation. Figure 9 shows a sample of the 3dshapes dataset. The number of values these attributes can take range from 4 to 14. We take 4 values from each attribute ($|\mathcal{X}| = 4$). For the attribute that correspond to the scale of the object, we choose values 0, 2, 4, 7 out of all the available values which range from 0 to 7. For the viewing orientation attribute, we choose values 0, 4, 9, 14 out of all the available values which range from 0 to 14. We construct each of the other attributes’ 4 values by random sampling.

Agent architecture The speaker processes the image with a two-layer convolutional neural network (CNN) (Schmidhuber, 2015) each of which is accompanied by a max pooling layer. The output then goes through a linear layer before being processed by the single-layer GRU as described in §4. This produces a message m . CNNs have kernel size of 8, stride of 1, and filter size of 8. We utilize same padding. Max pooling layer has kernel size of 2 and stride of 2. The linear layer projects activations of dimension 2,048 to 500. A listener with the same architecture as in §4 processes the message m and outputs its prediction for the values of the image’s attributes.

D Experimental details on ImageNet

The speaker processes the target image’s representation of dimension 2,048 with a linear layer producing activations of dimension 500. They are then processed by the single-layer GRU as described in §4. This produces message m containing descriptions of the target image.

A listener π_{ϕ_i} processes each of the images’ representations in its context \mathcal{C}_i with a linear layer then computes similarity scores of them with the message representation from the single-layer GRU described in §4. The message representation is computed from the last hidden state of the the single-layer GRU after it is passed through a lin-

ear layer. The resulting message representation has a dimension of 500. We use the dot product as the similarity score function. These scores are then passed to Softmax activation to produce a distribution over the context \mathcal{C}_i . We construct each listener’s context by randomly sampling images without replacement.

The vocabulary size $|\mathcal{W}|$ and message length T are both set to 10. The batch size is set to 2,048. Training is performed for 2,000 epochs and evaluation is performed with the checkpoint that exhibits the highest accuracy on the validation set. We repeat each experiment with 10 different random seeds and report the average. Scramble resistance is calculated with respect to one randomly selected listener. We report compositionality metrics from the test set. The image representations of ImageNet dataset are available under Apache-2.0 license.

E Detailed description of PosDis and BosDis

The entropy $\mathcal{H}(m_i)$ in the description of PosDis from §5 refers to the i -th position’s entropy in messages. Let $\mathbb{P}(m_i = w)$ denote the probability that i -th symbol of a message is equal to the symbol w . Then,

$$\mathcal{H}(m_i) = - \sum_{w \in \mathcal{W}} \mathbb{P}(m_i = w) \log \mathbb{P}(m_i = w)$$

where \mathcal{W} denotes the vocabulary.

Similarly, let n_i denote the number of occurrences of the i -th symbol in vocabulary \mathcal{W} . The entropy $\mathcal{H}(n_i)$ in the description of BosDis from §5 is written as

$$\mathcal{H}(n_i) = - \sum_{c=0}^T \mathbb{P}(n_i = c) \log \mathbb{P}(n_i = c)$$

where T is the length of messages and $\mathbb{P}(n_i = c)$ denotes the probability that the i -th symbol occurs c times in a message.

For both entropies, the probabilities are calculated by computing each symbol’s frequency in messages.

F Effects of relative model capacity in listeners of different interests regime

We validate that higher compositionality exhibited from listeners of different interests regime do not

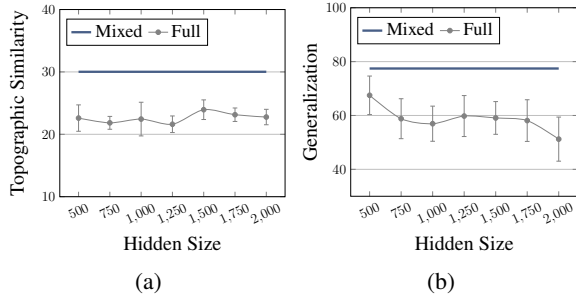


Figure 10: Language properties under varying values of listener hidden sizes in the full-interest formation in comparison with the mixed-interest formation of a fixed listener capacity.

stem from the relative difficulty of the task as the number of attributes that need to be determined is lower in that regime. To that end, we increase the hidden size of listeners in the full-interest formation from 500 to larger values and compare them with the mixed-interest formation with $K_i = 1$. The experimental setup follows from §4. The hidden size of listeners in the mixed-interest formation is fixed to 500. Both formations contain the same number of listeners, $N = 5$.

In Figure 10a, we observe that the values of TopSim stay almost the same as the listeners’ capacity is increased in the full-interest formation. This suggests that the relative capacity of the listeners in listeners of different interests regime is not the core contributing factor for the emergence of compositionality. Similarly, in Figure 10b, we observe a decrease in generalization ability as the capacity of the listeners in the full-interest formation is increased. These observations confirm that it is not the relative easiness of the task that induced more compositional languages in the listeners of different interests regime.

G Experiments on larger configurations

We test whether our findings hold on a larger number of attributes as well as vocabulary size. We double the number of attributes from 4 to 8 and the vocabulary size from 10 to 20. To keep the numbers in a manageable scale, the number of values an attribute can take is decreased from 10 to 5. This setup represents about a 39-fold increase in terms of the number of possible attribute-value combinations and a 32-fold increase in terms of the number of possible messages.

We conduct the same experiments as in §6.2 to see how listeners of different interests affect

language properties in this increased setup. We construct three kinds of listener formations. The single-interest formation contains eight listeners that are interested in each of the eight attributes ($K_i = 1$). The mixed-interest formation is the same as the single-interest formation but contains one additional listener that is interested in all of the eight attributes. The full-interest formation contains nine listeners that are interested in all of the eight attributes ($K_i = K$).

Figure 11 presents the results. The induced languages’ compositionality in terms of TopSim is higher when messages are more frequently read by listeners that are interested in smaller parts of the whole attributes. They also exhibit more symbol-wise structures rather than position-wise structures. Interestingly, the generalization ability of the mixed-interest formation is lower. We hypothesize this could be due to the larger ratio of single-interest listeners to the full-interest listener (8:1).

H Evaluation with concept-best-matching

Let \mathcal{C} be the set of possible concepts, where a concept represents a value manifestation in an attribute, so that $|\mathcal{C}| = |\mathcal{X}| \cdot K$. A weighted graph \mathcal{P} is constructed between the set of symbols, \mathcal{W} , and the set of concepts, \mathcal{C} , characterized by edges \mathcal{E} , where an edge $e_{i,j} \in \mathcal{E}$ represents connection from the i -th concept to the j -th symbol with a weight that is equal to the number of joint occurrences of the concept, symbol pair. An optimal pairing of symbols and concepts is sought, such that no two edges connect the same symbols or concepts, while maximizing the sum of edge weights. Concept-best-matching (Carmeli et al., 2024) metric is defined as a normalized sum of the edge weights of such an optimal bipartite graph induced from the full weighted graph \mathcal{P} . High concept-best-matching scores indicate unique association of a symbol to a concept.

We conduct evaluation with concept-best-matching metric on 3dshapes dataset and present the results in Table 3. The experimental setup follows from §7.1. Association of a concept with a symbol gets pronounced when listeners are interested in smaller subsets of the whole attributes, as can be seen by the larger values in concept-best-match score when the number of interested attributes is low. This trend is similar to the bag-of-symbols disentanglement results in Figure 6c.

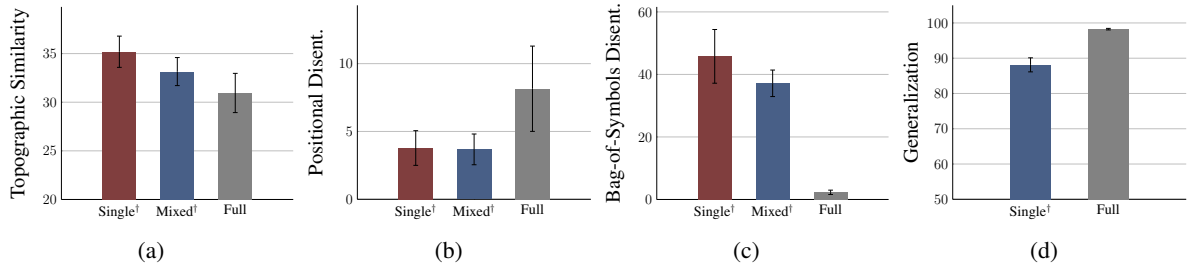


Figure 11: Effects of listeners of different interests in a larger configuration. [†] denotes statistically significant difference from the full-interest formation (one-tailed t -test with $p < 0.05$).

Listener Formation	Number of Interested Attributes				
	1	2	3	4	5
Partial-Interest	47.46±2.3 [†]	45.20±1.6 [†]	43.86±2.1	44.57±2.6	44.84±2.6
Full-Interest	44.86±3.1	43.88±2.2	45.01±2.0	43.88±2.2	44.86±3.1

Table 3: Evaluation results on concept-best-matching. [†] denotes statistically significant difference compared to the full-interest formation (one-tailed t -test with $p < 0.05$).

I Further discussion on limitations

Causes and implications of different compositionality structures In §6.2, we observe that listeners of different interests induce more symbol-wise structures in languages rather than position-wise structures, and we find a reverse trend when coordination pressure is exerted to the environment. We hypothesize that the effects may stem from the fact that an object cannot be described with the number of occurrences of a symbol when the number of attributes is large. This forces the language to describe an object with multiple symbols; hence the order in which the symbols appear may have a relatively higher chance of playing an important role. However, when a listener is only interested in a small subset of the whole attributes, these subsets may be sufficiently described by the number of occurrences of a symbol. Combined with the fact that GRUs do not utilize explicit positional encodings, may give rise to the preference of describing an attribute with occurrences of a specific symbol. A further investigation into the underlying mechanisms that cause these phenomena and their implications still remains to be conducted. Future work may also explore how these kinds of compositionality structures affect performance in downstream tasks from the perspective of representation learning.

Relationship to other environmental pressures

As we discuss in §2, there are various environmental factors involved in emergence of compositionality, e.g., noisy channel (Kuciński et al., 2021).

The relationship between these and the pressures investigated in this work could be further explored. For instance, we explore coordination pressure in relation to iterated learning in §6.4.

Effects of one-to-many communication on other language universals Our work focuses on one-to-many communication’s effects on compositionality. However, there are other language universals that are actively studied in the emergent communication field as discussed in §2. Future work may explore how one-to-many communication affects other language universals.

Availability of attribute labels In the experiments with listeners of different interests, the listeners’ interests are derived from labeled attributes. However, a dataset in question may lack such labels. Future work may investigate the ways in which interests can be formed in an unsupervised manner. One could devise information bottlenecks so that each listener would have a specialized role in the cooperative task.

J Reproducibility

For training, we utilized NVIDIA A100 80GB, NVIDIA RTX A6000 Ada, NVIDIA RTX A6000, NVIDIA RTX A4000, NVIDIA GeForce RTX 4090 and NVIDIA GeForce RTX 3090. The most demanding task in terms of compute required less than 24GB of VRAM and took about 2 or 3 hours to complete per random seed. The number of parameters of an agent is far less than 1M in all experiments.

We make our code publicly available at:
<https://github.com/hy18284/onetomany>.