



Universiteit
Leiden

The Netherlands

Collaborative meaning-making: the emergence of novel languages in humans, machines, and human-machine interactions

Kouwenhoven, T.

Citation

Kouwenhoven, T. (2025, October 30). *Collaborative meaning-making: the emergence of novel languages in humans, machines, and human-machine interactions*. SIKS Dissertation Series. Retrieved from <https://hdl.handle.net/1887/4281976>

Version: Publisher's Version

License: [Licence agreement concerning inclusion of doctoral thesis in the Institutional Repository of the University of Leiden](#)

Downloaded from: <https://hdl.handle.net/1887/4281976>

Note: To cite this publication please use the final published version (if applicable).

6

Searching for Structure

Human languages have evolved to be structured through repeated language learning and use. These processes introduce biases that operate during language acquisition and shape linguistic systems toward communicative efficiency. In this chapter, we investigate whether the same happens if artificial languages are optimised for the implicit biases of Large Language Models (LLMs). To this end, we simulate a classical referential game in which LLMs learn and use artificial languages. Our results show that initially unstructured holistic languages are indeed shaped to have some structural properties that allow two LLM agents to communicate successfully. Similar to observations in human experiments, generational transmission increases the learnability of languages, but can at the same time result in non-humanlike degenerate vocabularies. Taken together, this work extends experimental findings, shows that LLMs can be used as tools in simulations of language evolution, and opens possibilities for future human-machine experiments in this field.

Originally published as: Tom Kouwenhoven, Max Peeperkorn, Tessa Verhoef. 2025. Searching for Structure: Investigating Emergent Communication with Large Language Models. In Rambow, O., Wanner, L., Apidianaki, M., Al-Khalifa, H., Di Eugenio, B., Schockaert, S., editors, *In Proceedings of the 31st International Conference on Computational Linguistics*, pages 9977–9991, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

6.1 Introduction

Vocabularies of signals enable us to communicate about meanings, but to express an arbitrary number of meanings, vocabularies would require an equally large set of words as there are meanings, and learning such holistic vocabularies is cognitively challenging. Human languages, therefore, typically show some form of compositional structure, where meaningful signal-meaning mappings can be composed such that the combination of individual meaningful signals can express more than the meaning of the individual components alone (Hockett, 1960). An important finding in the field of language evolution is that such structural properties can emerge at the population level as a result of individual learning biases and pressures that continuously shape the languages on a longer timescale, often eventually resulting in languages that are easier to learn and exhibit some degree of structure (Smith, 2022).

The processes involved in the evolution of language have been extensively investigated through experiments and simulations. The latter typically use hard-coded agents with inductive biases (de Boer, 2006), Bayesian learners (e.g. Griffiths and Kalish, 2007b; Culbertson and Smolensky, 2012; Kirby et al., 2015), or reinforcement learning agents (Lazaridou and Baroni, 2020) to investigate the evolution of structured languages. In contrast, we investigate whether more flexible LLMs as relatively unbiased language learners (Wilcox et al., 2023) are appropriate tools to study how languages evolve. While their internal mechanisms are fundamentally different from those of humans, they still are the first close flexible comparators of human language users, which can be used as tools to answer cognitive and typological investigations (Warstadt and Bowman, 2022; van Dijk et al., 2023a). Given that languages are shaped by the biases and pressures of individual language learners, which differ for LLMs (e.g., fewer memory constraints), we are interested in identifying similarities and differences between humans and LLMs on specific language evolution-oriented tasks.

Our work largely follows the experimental design by Kirby et al. (2015) in which Bayesian learners and humans learn an artificial language to communicate in a referential game. They find that linguistic structure arises from a trade-off between pressures for compressibility and expressivity. This chapter extends their work by using LLMs as objects of investigation. Specifically, we investigate how artificial languages evolve when two LLMs communicate in a referential game and what the effects of generational transmission on these languages are. We compare the properties of these languages to those that are found in experiments involving humans. Results show that 1) LLMs can learn artificial languages and use them to communicate successfully, 2) the languages exhibit higher degrees of structure after multiple communication rounds, 3) LLMs generalise in more systematic ways when the evolved language is more structured, and 4) languages adapt, although not necessarily in a human-like way, and become easier to learn by the LLMs as a result of generational transmission.

6.2 Background & Related work

6.2.1 The evolution of structure

Learning novel signal-meaning mappings, and the emergence of rules that can combine these signals into structured languages have been abundantly investigated in the field of language evolution using human experiments (Kirby et al., 2008; Galantucci, 2005; Scott-Phillips et al., 2009; Verhoef, 2012; Raviv et al., 2019a,b) and computational simulations (de Boer, 2006; Steels and Loetzsch, 2012; Lazaridou and Baroni, 2020). These typically follow a setup where success depends on cooperation between two or more participants/agents in a Lewis game. Here, players are prevented from communicating using conventional communicative means and instead must establish novel communication systems through repeated cooperation. Outcomes often show that players, human or machine, quickly establish novel signal-meaning mappings that enable them to communicate successfully. However, recent computational simulations using reinforcement learning agents often develop communicative systems different from those of humans (Galke et al., 2022)¹ unless specific key pressures are introduced to recover initially absent human patterns (Galke and Raviv, 2025).

It has been suggested that seemingly arbitrary aspects of linguistic structure may result from general learning and processing biases deriving from the structure of thought processes, perceptuo-motor factors, cognitive limitations, and pragmatics (Christiansen and Chater, 2008). A well-investigated cause for this phenomenon is the process of cumulative cultural evolution (Boyd et al., 1996; Tomasello, 1999), which is typically investigated using iterated learning experiments (Kirby et al., 2008). Here, information (e.g., a language) is repeatedly passed down from one generation to the next, where the information is modified and improved upon within each generation. The influential work by Kirby et al. (2008, 2015) demonstrated that when human individuals learned an artificial language previously learned by another individual, the language became easier to learn and displayed a higher degree of structure. Crucially, these results are mostly attributed to the fact that the language repeatedly goes through a learning bottleneck, in which individual cognitive constraints, such as memory constraints, gradually shape the language. Iterated learning has been used to demonstrate that structure emerges in various setups with, for example, continuous signals (Verhoef, 2012) or continuous meaning spaces (Carr et al., 2017), and it is argued that this process may have led to the statistical Zipfian structure of language (Arnon and Kirby, 2024). Yet, Raviv et al. (2019a) showed that structure can also emerge *without* generational transmission. In this case, a pressure for compressibility originating from communication with multiple interaction partners and expanding meaning spaces causes languages to become compositional. This effect is even more prominent if the

¹But see Lian et al. (2023b, 2024); Zhang et al. (2024b) for recent work showing that the need to be understood (i.e. communicative success), noise, context sensitivity, and incremental sentence processing help induce human-like patterns of dependency length minimisation in reinforcement learning agents.

number of interaction partners is larger (Raviv et al., 2019b). The current chapter is inspired by the traditional methods described previously and extends them with our current most sophisticated models of natural language.

6.2.2 LLMs as models of language

LLMs are sophisticated models of natural languages, and growing evidence shows their ability to exhibit ‘average’ human behaviours. It is, for example, suggested that LLMs can model human moral judgements (Dillion et al., 2023) and transmission chain experiments revealed human-like content biases in GPT-3.5 (Acerbi and Stubbersfield, 2023). When LLMs are extended with records of experiences, Park et al. (2023) showed that groups of generative agents exhibit believable human-like individual and emergent social behaviours when they interact over extended periods. It is even suggested that human-LLM interactions in everyday life can potentially mediate human cultures through their influence on cultural evolutionary processes of variation, transmission and selection (Brinkmann et al., 2023; Yiu et al., 2024).

While previous work has investigated human-like behaviour at inference time, findings from cognitive science can also be used to improve model performance. Iterated learning can, for example, be incorporated into the training regime to extrapolate desirable behaviours. Zheng et al. (2024) have likewise shown that representations are easier to learn when vision-language contrastive learning is reframed as the Lewis signalling game between a vision agent and a language agent, ultimately improving compositional reasoning in vision-language models. However, this does not guarantee model improvements. Shumailov et al. (2024) have shown that LLMs, autoencoders and Gaussian mixture models drift when trained repeatedly on AI-generated data. In these cases, crucially, the generated content is slowly optimised to be understandable for models, *not* for humans, resulting in what they call model collapse. The authors therefore argue that genuine human interactions with systems will be increasingly important to prevent model collapse. While drift is often seen as an unwanted effect of unsupervised training, this is not surprising from a language evolution viewpoint since languages adapt to how they are learned and used (Smith, 2022). It was therefore suggested in Chapter 1 that languages should adapt to become more natural for humans *and* machines. This bears much resemblance to the idea that findings from cognitive science can prevent modal collapse (Smith et al., 2024) or inform modelling choices (Galke and Raviv, 2025). Here, we view LLMs from this evolutionary perspective.

Although inductive biases inherent to a language model’s (pre-)training objectives (i.e. the cloze task and instruction tuning) and memory constraints are very different from those in humans, recent work has shown that GPT-2 models struggle to learn languages that contain unnatural word orders, lack hierarchical structure, or lack information locality (Kallini et al., 2024). This suggests that, even though the language processing mechanisms in Transformers are non-humanlike, LLMs exhibit a preference for structured languages similar to those of

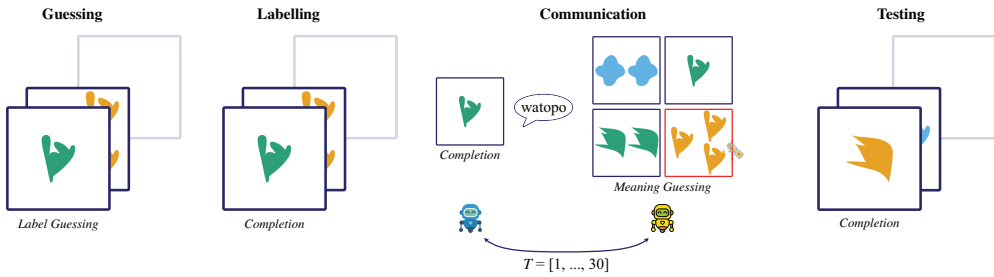


Figure 6.1: A graphical representation of the experimental blocks. The agents first go through a guessing block before labelling each of the 15 training stimuli in the labelling block. The communication block is done for 4 rounds, each consisting of 30 tasks T , where the agents alternate speaker-listener roles to be a speaker and listener for each stimulus once. Finally, the agents label 27 (15 original and 12 novel) stimuli in the testing block. Icons obtained from flaticon.com.

humans. Moreover, in an artificial language learning experiment similar to the work presented here, Galke et al. (2024) showed that compositional structure is advantageous for GPT-3 when learning an artificial language and that a higher degree of compositional structure also resulted in human-like generalisation for new unseen items. This chapter is different in that Galke et al. tested the ability of GPT-3 to learn languages that evolved during a *human* experiment (Raviv et al., 2019b, 2021), thus being optimised for human learners. We instead wish to investigate what kinds of languages evolve when they are optimised for *LLMs*.

6.3 Methodology

Our methodology is inspired by Kirby et al. (2015) and (Raviv et al., 2021). The complete simulation set-up consists of four blocks: guessing, labelling, communication and testing (Section 6.3.2 & Figure 6.1²). The agents perform the guessing, labelling, and testing block separately, but the communication block is interactive. The communication block is a classic referential game in which two agents communicate to discriminate a target stimulus from four distractor stimuli. They do so in four rounds, each consisting of 30 interactions T , alternating speaker-listener roles between interactions. In a single interaction round, the speaker observes a target stimulus (not the distractors) and utters a signal that describes the current stimulus. Using this signal, the listener must discriminate the correct target amongst a set of distractor stimuli. Cooperation is successful when the listener’s guess is the target stimulus.³

²This is for illustration purposes only, we stress that our simulations are entirely run in the textual modality only to avoid the additional challenge of extracting relevant visual features and mapping these to artificial languages.

³All code, materials, and data are available on OSF: <https://osf.io/52yar/>.

6.3.1 Stimuli and initial languages

The meaning space consists of stimuli with three attributes. They have one of three shapes, one of three colours, and can appear in groups of one, two, or three shapes, creating 27 distinct stimuli. Initial signals for these stimuli were generated before each experiment according to the method used by Kirby et al. (2008). The signals are concatenations of 2, 3, or 4 randomly selected consonant-vowel (CV) syllables resulting in artificial non-existing signals (e.g., watopo, nafa, nomomeme). The CV syllables consist of one of eight consonants g, h, k, l, m, n, p, w and one of five vowels a, e, i, o, u. Out of 27 stimuli, only 15 stimuli are used during the guessing, labelling, and communication blocks. All 27 stimuli are used in the testing block such that we can assess whether the agents can generalise to novel stimuli. The training stimuli are selected randomly before each simulation, but we ensure that each attribute value is represented equally often across this set.

6.3.2 Simulation blocks

Each simulation consists of four blocks. In the first block, we assess whether agents can correctly guess a signal when presented with a stimulus. Second, in the labelling block, an agent repeatedly produces a signal for each stimulus given the initial training vocabulary. The signals generated in this block are taken as the learned vocabulary for that agent. In the third block, the agents communicate as described before, taking turns as speaker and listener until all rounds are completed and each stimulus appears twice per round (i.e., both agents produce a signal for each stimulus and make a guess for each stimulus). In this block, the interaction between the agents gradually alters each agent's individual vocabulary, much like how this is done in earlier simulations (De Boer, 2000; Steels and Loetzsch, 2012). Specifically, we update the current stimulus to be associated with the signal that is produced. After the communication block, the testing block tasks the agents to generate signals for the entire meaning space of 27 stimuli using the training vocabulary that was optimised in the labelling and communication block. Hence, they must generalise their strategies to unseen samples.

6.3.3 LLMs as agents

The LLMs used in our experiment were instruction-tuned instantiations of Llama-3 70B (Llama Team, 2024) with greedy sampling.⁴ Since our method required LLM agents to follow instructions, we did not consider base models. In particular, we instructed them about the nature of their task and its collaborative goal. Though instruction-tuning using reinforcement learning from human feedback (RLHF) may influence the probabilities of some tokens fitting to instruction-following behaviour, the capacity to produce fluent language and knowledge

⁴Although we only report results on one model type, initial explorations with GPT-3.5 and Llama 2 7B showed similar behaviours to Llama-3 70B.

```
{ 'shape': 3, 'colour': 'blue', 'amount': 1, 'word': 'ninikonu' }
{ 'shape': 1, 'colour': 'green', 'amount': 3, 'word': 'hanosa' }
  ⋮
{ 'shape': 2, 'colour': 'orange', 'amount': 2, 'word': 'sanu' }
{ 'shape': 1, 'colour': 'green', 'amount': 3, 'word': '[COMPLETE]
```

Prompt 6.3.1: A vocabulary snippet as used in a completion prompt. The complete prompts are visible in Section C.2.

is mostly acquired in the pre-training phase (Zhou et al., 2023; Lin et al., 2024). Moreover, since our method does not specifically tap into instruction-tuning behaviour, we do not expect much variance in the results should we use base models only. While human participants typically learn signal-meaning mappings through a learning block, we use LLMs’ in-context learning ability (Brown et al., 2020) to teach them the languages. Specifically, we prepend our prompts with the items to be learned in a structured JSON-like format (Prompt 6.3.1). Given the observed behavioural similarities between humans and LLMs (Galke et al., 2024), we assume that a vocabulary of signal-meaning mappings in the context of a prompt provides enough (distributional) information for a LLM to learn an appropriate mapping between the attributes of the stimuli and signal syllables. Although the prompt structure ‘invites’ the LLM to infer a signal from the stimulus attributes, we are agnostic about how exactly and what kind of mapping the LLM deduces, but we are interested in the resulting behaviours.

Throughout a simulation, agents essentially perform one of two tasks: generation or guessing. The labelling block and speaking in the communication block involve generating signals. The guessing block and discrimination in the communication block involve guessing. The prompts for these tasks are extensions of those used by Galke et al. (2024), with slight adaptations to enable LLMs to discriminate between stimuli. Given that LLMs show a primacy and recency bias (Liu et al., 2024), the vocabulary is shuffled before each task such that ordering effects are minimal. System instructions depend on the task performed, but are largely similar and chosen to be as close as possible to instructions given to humans in experimental settings.

Generating signals. For signal generation in the labelling block, we use prompt completion (Prompt C.2.1). During labelling, the agents see the *entire* training set and generate a signal for each stimulus, effectively amounting to a look-up task since the stimulus is present in the prompt. On the other hand, the vocabulary presented to agents during communication and testing does *not* include the current stimulus, thus requiring the agents to extract an appropriate mapping and generalise to new stimuli (Prompt C.2.2). A human-like solution would be to map stimulus attributes (i.e. shape, colour, and amount) to syllables representing these attributes and create compositions that describe the stimulus. During communication, we incentivise the agents to communicate using a `communicativeSuccess` attribute which is set to 1 if the previous interaction for this stimulus was successful and zero otherwise. Adding this

attribute functions as a memory between interactions and provides a pressure for expressivity. It is hypothesised that the latter plays an important role in human language evolution since it prevents languages from becoming degenerate (Smith et al., 2013). Importantly, during testing, the vocabulary presented to the agents always includes the stimuli present in the train set (without the current stimulus), and stimuli from the test set are never present.

Guessing signals or meanings. For guessing and discrimination during communication, the agents need to respond with a choice corresponding to the speaker’s signal. Unfortunately, LLMs are inconsistent and unreliable in answering multiple-choice questions (Khatun and Brown, 2024). In our initial exploration, this indeed proved to be unusable. Instead, for each distractor (signal or meaning), we run the prompt prefilled with that distractor through the model and select the distractor with the highest probability (Prompt C.2.3). Again, the agents observe the training vocabulary *with* the current stimulus in the guessing block. In the communication block, agents observe the training vocabulary *without* the current stimulus.

6.3.4 Metrics

6

We are firstly interested in investigating whether two agents settle on a language that enables them to communicate, measured by the percentage of successful interactions (*PercCom*) in a round. We use multiple metrics to measure structure in messages. The most common metric is topographic similarity (*TopSim*, Brighton and Kirby, 2006). Similar to Kirby et al. (2008), we report Z-scores of the Mantel test (Mantel, 1967) between signal similarities (normalised Levenshtein distance) and semantic similarities (the number of equal attributes between two meanings). A communication system with a high *TopSim* uses similar signals for similar meanings. We compute the *Ngram* diversity (Meister et al., 2023), being the average fraction of unique vs. total *Ngrams* for $N \in \{1, 2, 3, 4, 5\}$ in all produced signals. Low *Ngram* diversity across all signals implies the agents re-use parts of signals in different signals, hinting at compositional signals when it happens in combination with increased *TopSim*. We assess the degree of signal systematicity between the signals produced for unseen stimuli in the test block and the previous stimuli in the communication block using the generalisation score (*GenScore*, Raviv et al., 2021). Here, we first compute the pairwise semantic difference between each stimulus in the train and test scenes, followed by the pairwise normalised edit distance between the signals produced for these scenes. We then take the Pearson correlation between these differences across all stimuli. Intuitively, this measures whether similar scenes across both sets are similarly labelled, thereby suggesting generalisation.

6.4 Evaluation

We ran 15 simulations, each initialised with a random seed and unique artificial, unstructured, and holistic language.⁵ Metrics were computed for each block, except for the generalisation score, which is only computed for the testing block. A human-like result would show increasingly successful interactions and increasing *TopSim* scores, while *Ngram* diversity should go down. If this is the case, we expect to observe higher generalisation scores since agents can compose new signals according to a learned structured strategy. We use linear mixed effects models to analyse the results of the communication block and to account for the random effects of each simulation's vocabulary. The slope ($\hat{\beta}$) determines the direction of the effect and the rate of change. Additionally, we use conditional R^2 (Nakagawa and Schielzeth, 2013), denoted by R_c^2 , which considers fixed and random effects, to show how much variance can be explained by the model. Higher values of R_c^2 indicate that the model captures more variance and that correlations are stronger. Finally, we report the marginal R_m^2 , which is the variance explained by the fixed effects.

6.5 Results

6

6.5.1 Learning the artificial languages

We first assess whether LLMs were able to learn the initially unstructured languages. Given the nature of the guessing task, which is essentially a lookup task, unsurprisingly, LLMs were able to guess the correct signals for the stimuli almost perfectly ($M = .973, SD = .031$). However, labelling the same stimuli via completion proved much more difficult ($M = .453, SD = .152$) despite the presence of the correct signal in the prompt. This contrast is in line with work showing that LLM predictions are sensitive to task instructions and how predictions are extracted (Weber et al., 2023; Hu and Levy, 2023; Hu and Frank, 2024). Additionally, it corroborates the use of prefilled options in our guessing prompts during communication. Nevertheless, this performance is still better than that of humans⁶ and is not unimpressive given the vast number of possible signals that can be produced. Finally, the expected struggle to correctly reproduce (i.e., learn) unstructured signals introduces some welcome variation to the agents' vocabulary, which is used at the start of the communication block.

⁵We are aware of the fragile nature of behavioural experiments with LLMs. Small perturbations to prompts can have large effects on the outcome (e.g. Weber et al., 2023; Hu and Levy, 2023; Hu and Frank, 2024; Giulianelli et al., 2024). This is also the case in our experiment. To ensure the reproducibility of the current findings, we use an open-source model, share all prompts, log probabilities, and data on OSF. Nonetheless, the probabilistic nature of LLMs will always warrant further investigation.

⁶In Chapter 7 we conduct an experiment involving humans and show that the guessing block is much easier than the labelling block.

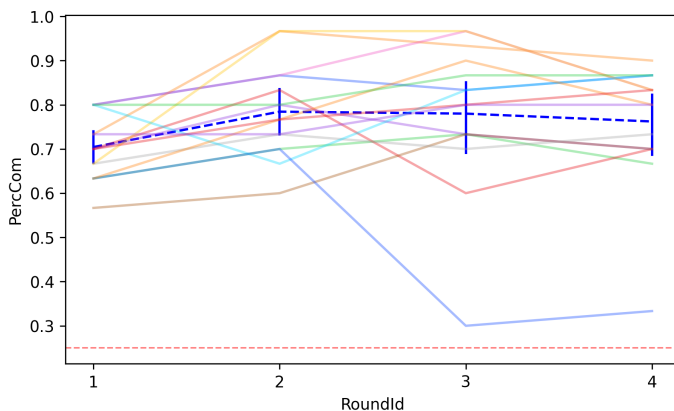


Figure 6.2: The communicative success (*PercCom*) over the communication rounds. Each coloured line indicates a simulation, and the dashed blue line displays the average with bars indicating the 95% confidence interval. The dashed red line delineates chance performance.

6.5.2 Agents communicate successfully

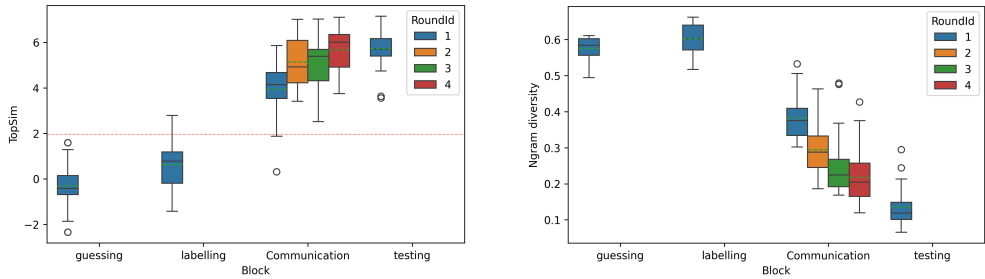
6

Once the agents have individually learned the vocabulary, they start communicating. Despite initially starting with different languages, approximately 70% of the interactions in the first round are successful (chance performance would amount to 25%). This increases somewhat in the following rounds to $\approx 75\%$, but not significantly (Figure 6.2). Interestingly, communicative success is not guaranteed; it fluctuates between rounds and can even decrease drastically in some simulations.

6.5.3 Communication results in structure

Although the initial languages are unstructured, some form of structure emerges due to repeated learning and use (Figure 6.3). This mostly happens during the communication block where *TopSim* increases significantly across rounds ($\hat{\beta} = .508 \pm .073$, $R_c^2 = .579$, $R_m^2 = .355$, $p < .001$) and *Ngram* decreases across rounds ($\hat{\beta} = -.054 \pm .004$, $R_c^2 = .812$, $R_m^2 = .558$, $p < .001$). This increase in structure benefits communicative success positively ($\hat{\beta} = .035 \pm .007$, $R_c^2 = .769$, $R_m^2 = .427$, $p < .001$). However, we also observe behaviour that is not human-like; the signals used to communicate become longer over the rounds ($\hat{\beta} = .557 \pm .044$, $R_c^2 = .919$, $R_m^2 = .505$, $p < .001$). This contradicts what is observed in human experiments, where we typically observe that messages become shorter and lie close to a theoretical frontier balancing expressivity and simplicity (Piantadosi et al., 2011; Kirby et al., 2015).

These results extend the findings of Galke et al. (2024) in that LLMs not only *learn* structured vocabularies better but also naturally *shape* languages to have some form of structure when they are optimised for their inherent preferences. In addition to the fact that LLMs struggle to learn



(a) TopSim scores over the agent’s vocabulary in each block and round. The dashed red line indicates the $p < .05$ level.

(b) Ngram diversity scores over the agent’s vocabulary in each block and round.

Figure 6.3: Communication clearly increases the structure of the vocabularies, as seen by the increasing *TopSim* scores and decreasing *Ngram* diversity.

impossible languages (Kallini et al., 2024), that reframing prompt instructions into a structured list improves the model response (Mishra et al., 2022), and given that we do not impose pressure to induce structure, the surprising outcome of our experiments may be the result of an apparent “structure bias” in LLMs.

6.5.4 Structure enables better generalisation

After the communication block, the agents engage in the final simulation block. Here, they generate signals for all 27 stimuli using the vocabulary that has evolved after learning and communication. We find that high *TopSim* languages allow for better generalisation ($r = 0.735, p < .001$, Figure 6.4). A qualitative inspection of the signals generated in the testing block of the simulation, which resulted in the highest *TopSim* after communication, reveals that this agent repeatedly re-uses parts of signals in different compositions (Table 6.1). For example: “su” refers to the amount one, “pepi” to two, “petite” to three. For shape 1, the signals “sunu” and “sutu” are used, “ginu” for shape 2, and shape 3 is referred to with “wipi” or “wipu”. However, colours are less clearly demarcated by unique signal parts. This is also reflected in the ratio of unique signals produced during the test block ($M = 62.1\%$, $SD = 19.8\%$), indicating that some simulations sometimes result in repetitive use of the same signals for different meanings, resulting in a somewhat degenerate vocabulary. Such ambiguity may be the reason for non-perfect communicative success during communication. Nevertheless, it is clear that unseen stimuli are often labelled similarly to previously seen stimuli.

	Shape	Colour	Amount	Word
train set	3	orange	1	wipisu
	1	green	2	sutupepi
	2	green	1	ginisu
	3	green	1	wipisu
	1	blue	2	sunupepi
	1	green	3	sutupitite
	2	orange	1	ginusu
	3	blue	3	wipipitite
	3	green	3	wipupitite
	3	blue	1	wipisu
	1	blue	3	sunupitite
	2	orange	3	ginupitite
	2	blue	2	ginupepi
	1	orange	2	sunupepi
	2	orange	2	ginupepi
test set	1	orange	1	sutisu
	1	orange	3	sutupitite
	1	green	1	sutusu
	1	blue	1	sunusi
	2	green	2	ginupepi
	2	green	3	ginupitite
	2	blue	1	ginisu
	2	blue	3	ginupitite
	3	orange	2	wipupepi
	3	orange	3	wipipitite
	3	green	2	wipupepi
	3	blue	2	wipupepi

Table 6.1: An exemplary vocabulary that evolved in a simulation where the signals produced in the testing phase resulted in the highest *TopSim* score (7.13) after communication. The signals for the test stimuli share parts of signals and are composed similarly to train stimuli (*GenScore* = .792).

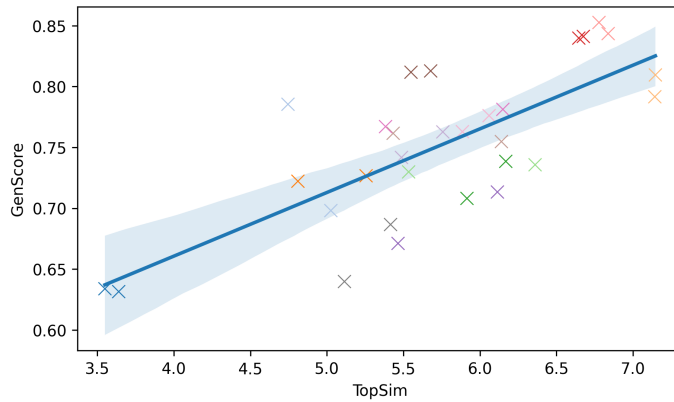


Figure 6.4: Languages that have evolved to be more structured allow for better generalisation to unseen test stimuli. Coloured crosses refer to individual simulations.

6.6 Iterated learning

The previous results showed that two LLMs can successfully communicate and slowly shape the language to become more structured. Provided that cumulative cultural evolution can extrapolate weak biases to have strong effects in socially learned systems like language (Smith, 2011), we extend our simulations by adding generations of learners. The first generation is initialised with a random unstructured language described in Section 6.3.1, but in the following generations, agents learn a portion of the signal-meaning mappings produced in the testing block by the agents of the previous generation. Only the vocabulary of the agent with the highest *TopSim* is transmitted to the next generation. We ran six transmission chains, each consisting of 8 generations. The seed generations for each chain were selected randomly from our initial 15 simulations.

6.6.1 Learnability increases

Iterated learning clearly increases the learnability of vocabularies (Figure 6.5). While LLMs in the first generation struggle to look up signals and reproduce them, a single generation of learning and using a language tremendously decreases the edit distance between ground truth signals and the produced signals. These results are remarkably similar to findings with human participants (Kirby et al., 2015), and show that the languages are optimised for LLMs' preferences.

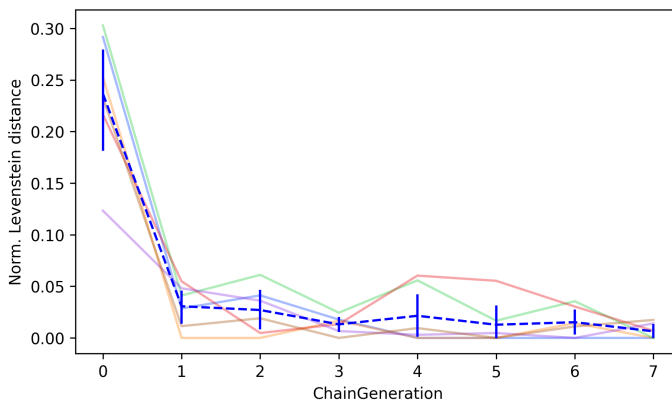


Figure 6.5: The normalised Levenshtein distance between the ground truth and the produced signal in the learning block. Solid lines indicate chains, and the dashed blue line indicates the average Levenshtein distance across simulations in a generation.

6.6.2 Communicative success and non-humanlike structures

6

Despite the increase in learnability, we do not observe an increase in communicative success due to iterated learning (Table 6.2, Figure C.1). This is possibly due to the already high scores of the first generation. Despite their increased learnability, the signals become significantly longer and more ambiguous. We take this non-humanlike solution to be an artefact of an absence of pressures for memorisation in LLMs. While human language is optimised to be compressible and expressive (Fedzechkina et al., 2012; Tamariz and Kirby, 2015; Kirby et al., 2015), the context windows of LLMs are considerably larger. In our case, Llama-3 70B has a context window of 8.2K tokens, which we do not exceed and therefore does not induce a pressure for compressibility.

Finally, the metrics to measure structure display a mixed picture. *TopSim*, does increase across generations but not significantly (Table 6.2). Yet, *Ngram* diversity decreases significantly across generations. For the evolution of these metrics across generations, see Section C.1. Qualitative inspections of several vocabularies show that some languages evolve into degenerate languages with repeating signals for different stimuli (i.e., underspecification). This is corroborated by a significantly lower number of uniquely produced signals in the last generation compared to the first simulation ($t(5) = 2.64, p = .046, M_{gen0} = .707, SD_{gen0} = .142, M_{gen7} = .519, SD_{gen7} = .119$). Together, this causes the *Ngram* diversity to be lower while clearly hurting communicative expressiveness. Even though degenerate languages are not uncommon in iterated learning experiments with humans (e.g., experiment 1 in Kirby et al., 2008), an additional pressure for expressivity typically prevents languages from becoming underspecified. Given the expressivity pressure that we imposed during the communication block, we expected to see less of such underspecification. The process of iterated learning, therefore, results in vocabularies

	$t(5)$	p	M_{gen0}	SD_{gen0}	M_{gen7}	SD_{gen7}
<i>PercCom</i>	.308	.770	.769	.077	.785	.123
<i>TopSim</i>	-1.42	.215	9.62	1.21	10.5	1.77
<i>Ngram</i>	2.83	.037	.158	.074	.071	.025

Table 6.2: The descriptives and statistics of the first (*gen0*) and last generation (*gen7*) in our chains. Paired t-tests show that *Ngram* diversity does significantly change resulting from generational transmission, while *TopSim* and *PercCom* do not.

that are optimised for the preferences of LLM agents but do so in a non-humanlike way.

6.7 Discussion

Our findings present a mixed picture; agents comprised of LLMs can learn and use artificial languages in a referential game. They do so by optimising the initially holistic vocabulary to fit better with the preferences of their language model, resulting in increased regularity and structure (Table 6.1). These human-like results are much in line with previous findings showing that structured languages can emerge from repeated interactions between interlocutors (i.a. Selten and Warglien, 2007; Verhoef et al., 2016b; Nölle et al., 2018; Raviv et al., 2019a). Yet, we also observe some degeneracy, i.e., many-to-one mappings of signals and attributes, and non-humanlike behaviours such as a tendency to produce long signals. Iterated learning further increases the learnability of the vocabulary but also extrapolates these non-humanlike behaviours further. Despite not being able to *directly* compare our results to human data, these findings are loosely comparable to earlier work involving human participants (Kirby et al., 2015; Raviv et al., 2019b) in which languages with similar properties emerge.

Table 6.1 moreover suggests that certain attributes, such as the colour attribute, in the inputs may be ignored, possibly due to the primacy and recency bias in LLMs (Liu et al., 2024). Optimising the instructive sentences by choosing sentences that maximise the fraction of valid model answers for each task, as suggested by Aher et al. (2023), may alleviate these ignorances and increase focus on relevant attributes. It is also possible that the LLMs do not ‘experience’ enough pressure to be understood by other agents, i.e., the *communicativeSuccess* attribute is not able to force a need to be expressive, which is deemed an essential pressure in computational simulations for human-like structures (Galke and Raviv, 2025). Despite these discrepancies, it is nevertheless interesting that some form of structure emerges.

Our results furthermore show variability between generations of learners. This is not uncommon in human experiments where processes of interaction and transmission sometimes generate fully systematic, compositional languages, but can also result in systems that lack structure entirely (Verhoef et al., 2022). In Chapter 2 we showed that differences in personal

biases may be a contributing factor to these differences. Since we do not initialise agents with different biases, these variations, originating in distributional information of the prepended vocabularies, are a natural human-like outcome of repeated exposure to and use of the language.

The evolution of degenerate vocabularies could be explained by the use of greedy decoding during signal generation, which does not necessarily produce the most human-like text (Holtzman et al., 2020; Meister et al., 2022, 2023) and may therefore also result in non-humanlike composition. Once an agent, perhaps mistakenly, duplicates a signal, its raw probabilities are increased when producing the next utterance, possibly resulting in a feedback loop that collapses onto a degenerate vocabulary. This effect may be further increased due to LLMs' inability to innovate (Bender et al., 2021; Yiu et al., 2024) and the choice of structured prompts that do not explicitly ask for innovation. Future work could attempt to increase the composition of novel signals by increasing the temperature parameter. Perhaps resulting in slightly more novel outputs as this forces exploration of the vocabulary embedding space (Peeperkorn et al., 2024), possibly alleviating the evolution of degenerate vocabularies and shifting the optimisation of the language to different solutions.

The rapid increase in learnability resulting from iterated learning proves that weak learning biases in language models, such as an observed simplicity bias (Chen et al., 2024), can be amplified by the process of generational transmission. Simulations with increased communicative difficulty, e.g., by increasing the number of distractors or the number of interaction partners, could reveal whether and how some form of memory constraint affects the learnability of languages, while also capturing the diversity and dynamic nature of language in the world more accurately. In general, systematic manipulations across model features (e.g., size, training data, or decoding strategies) may expose why we observe tendencies such as producing longer signals. Similar to what was proposed by Galke and Raviv (2025), we argue that careful manipulation of our setup can help reveal underlying mechanistic biases of language models and inform modelling choices when simulating language acquisition in LLMs. Taking into account the important role communication plays in shaping human language, LLM performance drastically increased when it was optimised for successful communication through RLHF.

Finally, we acknowledge that our results depend on several methodological considerations, including the model used, the prompt format, task instructions, and the tokenisation process. However, our primary goal was to investigate whether LLMs can be used in simulations of artificial language emergence. We aimed to stay as closely as possible to well-established experimental methods in the field of language emergence. We did not optimise for performance, human-like results, or compositional vocabularies. Instead, our goal was to reveal the natural behaviours of LLMs resulting from learning and using artificial languages. Future work could extend our findings by performing experiments in which humans collaborate with LLMs to investigate whether languages can evolve that are optimised for human *and* LLM preferences. Finally, as this chapter focused on experiments with a single LLM, future research should verify these findings across multiple LLM architectures to establish their generalisability.

6.8 Conclusion

Given the remarkable linguistic abilities of recent LLMs, we show how LLM-augmented agents behave in a classical referential game in which artificial languages, typically used in the field of language evolution, are learned and used. Primarily, our results suggest that LLMs can be used as artificial language learners to investigate the evolution of language. We showed that initially unstructured languages are optimised for improved learnability and allowed for successful communication. While we found some evidence of human-like compositional structures that enhance generalisation abilities, we also identified notable differences in the behavioural characteristics of LLMs compared to humans. Notably, iterated learning processes increased vocabulary learnability but also amplified such different characteristics further. As such, we extend existing research by revealing that structured languages are not merely easier for LLMs to learn. Critically, the inherent biases of LLMs also shape unstructured languages towards increased regularity. These findings contribute to a deeper understanding of how LLMs process and evolve language, potentially bridging the gap between computational models and natural language evolution. Finally, we hope to have shown that our setup is helpful in exposing the underlying mechanistic biases of LLMs and demystifying their uninterpretable nature.