
Compositional Communication with LLMs and Reasoning about Chemical Structures

Anonymous Authors¹

Abstract

Compositionality of communication is considered a prerequisite for reasoning. Despite overall impressive performance, LLMs seem to have fundamental issues with compositionality in reasoning tasks. Research of the emergence of languages in referential games demonstrates that compositionality can be achieved via combination of the game organization and constraints on communication protocols. In this contribution we propose and offer initial evaluation of the hypothesis that compositionality in reasoning tasks with LLMs can be improved by placing LLM agents in the referential games that coax compositionality of the communication. We describe a multi-stage chemical game including recognition, naming, and reconstruction of chemical structures by LLM agents without leveraging their pre-existing chemical knowledge.

1. Introduction

Reasoning is the hallmark of scientific process. Scientific applications of AI are yet to include seamless collaborative reasoning with human scientists. Specifically, compositionality appears to represent a big challenge even to the models with otherwise outstanding capabilities. Fundamentally, reasoning tasks might require some new AI architectures. Meanwhile, it pays to understand how much LLMs can be pushed before they reach a performance ceiling. Our approach is informed by the body of research of emergent communication in multi-agent reinforcement learning (MARL) (Lazaridou & Baroni, 2020). It is established that compositionality of the emergent languages is an independent feature that can be achieved via special constraints on the communication protocol and/or specific organization of the game where communication unfolds

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

(Chaabouni et al., 2020). We hypothesize, that LLMs communication can be pushed towards higher compositionality if LLMs are trained or fine-tuned as they participate in a properly organized referential game. Of course LLMs already have a handle on the natural human language and the game is not expected to produce a new language. The role of the game is to coax LLM agents to prioritize compositional communication over non-compositional (Kottur et al., 2017; Korbak et al., 2019).

We focus on reasoning tasks associated with chemical structures. Chemistry is one of the most successful applications of AI in science. Molecular organic chemistry is highly composable which makes it especially favorable for AI-driven applications. On top of it, the phrasing of theory of chemical structure is essentially identical to the definition of compositionality:

- Theory of chemical structure: the chemical nature of a complex molecule is determined by the nature and number of atoms, its components, their mutual arrangement and interaction.
- Compositionality: the meaning of any complex expression is determined by the meanings of its parts and the way they are put together.

LLMs struggle with composability of chemical structures and compositionality of reasoning about chemical structures. The issue is quite pressing because the majority of relevant chemical discovery workflows require a seamless, peer-like interaction of AI with human chemists along the lines of reasoning about components of molecular structure and their impact on utility of molecules.(Ristoski et al., 2020)

Referential games are often used to study emergence of languages (Lewis, 1969). Here, we are considering an asymmetric game with two agents, the Sender and the Receiver. As the Sender is exposed to the objects in the world, it learns to represent these objects and to associate utterances with the representations. The Sender shares utterances with the Receiver over a communication channel which in our case is discrete, variable length, and noiseless. The Receiver learns to associate utterances with its own representation of the world objects and to reconstruct the world objects. In

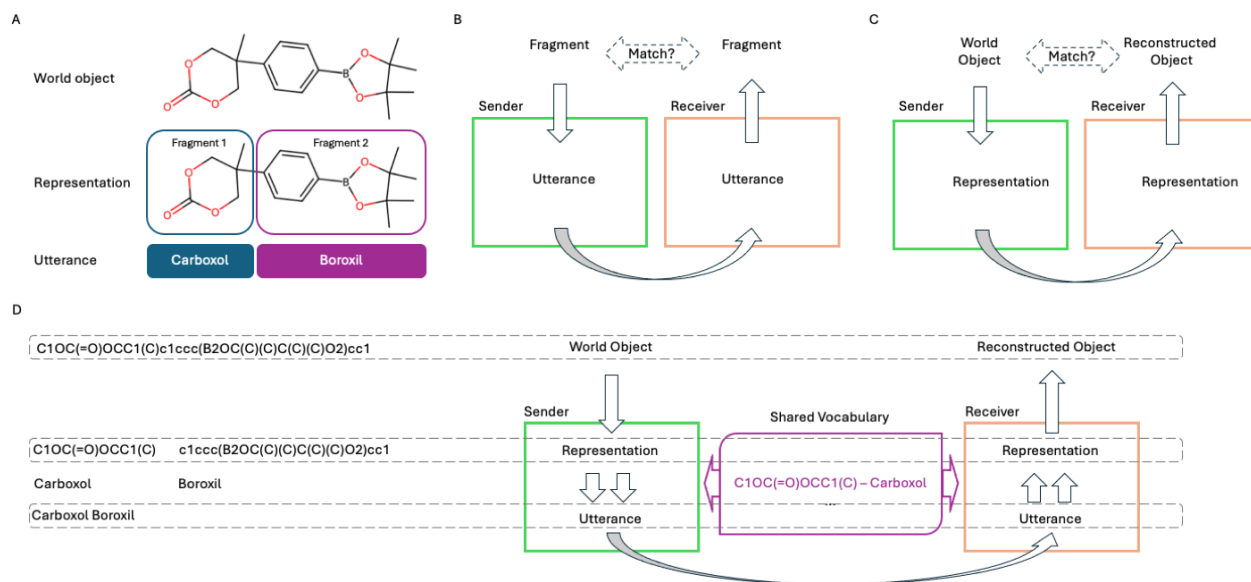


Figure 1. Complex referential games, such as template transfer, have been shown to support emergence of compositional communication [REF Template transfer] about multi-attribute objects. We are considering a similar construct, where objects, such as molecules, are composable from fragments, and the compositional referential game can be scaffolded on simpler games. **Panel A.** World objects are SMILES strings representing molecules. SMILES are split into substructure SMILES corresponding to the polymerizable groups (Fragment 1) and pendants (Fragment 2). Fragments are assigned names in the first sub-game where Sender and Receiver established a shared vocabulary. **Panel B.** First sub-game is a simple signaling game. The Sender and the Receiver learn a shared vocabulary by naming a finite library of molecular fragments. **Panel C.** Second sub-game involves decomposition of the composable object into fragments. The agents learn object representation as concatenated sub-strings corresponding to the fragments from the shared vocabulary. **Panel D.** Final nested referential chemical game with shared vocabulary where the Sender decomposes a composable object into fragments from the shared vocabulary, retrieves the names of these fragments, and constructs the utterance for the Receiver and the Receiver goes through the inverse process.

MARL settings, the agents are rewarded for each instance of communication where the Receiver correctly identified the object that the Sender was exposed to. Accumulated rewards are used to train the agents to improve their performance.

In this contribution, we train LLM model via fine-tuning on the pairs object-representation, representation-utterance, utterance-representation, and representation-object. The general structure of our chemical referential game closely follows (Korbak et al., 2019) and, by extension (Barrett & Skyrms, 2017). The world objects are SMILES strings that can be seen as simple concatenation of SMILES substrings. They are constructed as a combinatorial library from a set of performance-relevant groups and pendant groups. Each SMILES in the world is described with a two-part message, where the first part designates the performance-relevant group and the second part designates the pendant. This choice is informed by the referential games handling multi-attribute objects, cf. shape-color in (Korbak et al., 2019).

1.1. Related work

Our effort exists at the intersection of three active areas of research: reasoning and compositional communication with

LLMs, emergence of compositional languages in MARL, and application of LLMs in chemistry. It's been demonstrated that while most invented languages are effective (i.e. achieve near-perfect rewards), they are decidedly not interpretable or compositional (Kottur et al., 2017). This study showed development of the compositionality as a response to limiting vocabulary and eliminating memory of one of the communicating agents to simplify development of shared grounding of communication tokens. Another study (Korbak et al., 2019) reported achievement of emergent compositional communication in a complex game organized according to a template transfer pattern introduced in (Barrett & Skyrms, 2017). Elicitation of compositional generalization capabilities from LLMs has been approached via development of prompting strategies, such as skills-in-context (SKiC) prompting (Chen et al., 2023), which instructs LLMs how to compose basic skills to resolve more complex problems. A novel prompt-free approach, Compositional Task Representations (CTR) (SHAO et al., 2023), views each task as a composition of latent codes. CTR employs multi-task training to learn a discrete, compositional codebook, substantially outperforms prompt-based methods in zero-label learning on average.

Introduction of chemical benchmarks for LLMs ((Guo et al., 2023)) revealed general difficulties in comprehension of SMILES notation which translates into issues in downstream tasks. Focus of chemical applications of LLMs on instructions inevitably runs in the bottleneck of handling composability and compositionality of chemistry.(Fang et al., 2024)

2. Methodology

2.1. Data

Molecular combinatorial library is constructed from 7 performance-relevant groups and 63 pendant groups. The groups are merged via simple concatenation either in "performance group" + "pendant group A" pattern or "performance group" + "pendant group A" + "pendant group B" pattern, producing total of 11042 SMILES strings suitable for LLM fine-tuning. Only the first pattern including two fragments per molecule is used in the referential game setting 1A mimicking the games involving objects with two attributes (Korbak et al., 2019; Lazaridou & Baroni, 2020).

2.2. Game

The first sub-game 1B is a simple signaling game where the Sender and the Receiver establish a shared vocabulary that allows them to communicate about a fixed set of objects, such as fragments that are used to construct the combinatorial library. In the studies of language emergence, the agents are free to converge on any arbitrary vocabulary. In our case, both LLM agents are exposed to the natural language, scientific terminology and even SMILES notation. It is tempting to leverage SMILES comprehension to name the fragments. However, inconsistent handling of composability of SMILES is the reason why we are setting up the game-based process in the first place. Therefore, we proceed by asking the Sender to come up with short, unique names for the fragments that are not established chemical terms. The Receiver then needs to learn the correspondence between names and fragments. Effectively, the Receiver faces a supervised learning task on a small dataset, so for practical considerations we simply included the look-up table of fragments and names in the system prompts of both LLM agents and instructed the agents to use the table for search and retrieval of the relevant items.

In the second sub-game 1C the Sender learns to split a SMILES string into the sub-strings in the shared vocabulary. This primary task implies the secondary task, where the Sender has to match the fragment strings produced during the split to the content of the look-up table in the system prompt, and if both fragments have exactly matching entries, the Sender has to retrieve the corresponding names from the table. The Receiver handles the similar inverse task,

except that it needs to split a space-separated name shared by the Sender instead of a single SMILES string which is an enormous simplification.

These sub-games are nested in complete referential game 1D. The Sender encounters a world object, represents it as a set of fragments that have exact matches in the shared vocabulary, retrieves names of these fragments, and combines the names into a message. The Receiver parses the message into names of the fragments, retrieves the fragments from the look-up table, and reconstructs the world object.

2.3. Model training

The language model used as sender and receiver in this signaling game was fine-tuned on a dataset derived from data described in section 2.1. From the 11,042 SMILES strings and associated performance + pendant group labels in the Molecular combinatorial library, we created a dataset of input and output texts. This dataset covers various tasks that help LLMs learn to: a) split an initial SMILES notation of a molecule into sub-structure SMILES, b) map sub-structure SMILES to fragment names, c) map fragment names to sub-structure SMILES, and d) construct a SMILES string from the sub-structure SMILES of its fragments. We used Meta-Llama-3-70B-Instruct (AI@Meta, 2024) to create prompt variations for all four tasks, resulting in a dataset of 103,300 entries for fine-tuning the LLMs.

This work utilizes two different LLMs: 1) Phi-1.5 (Li et al., 2023), a small-sized model with 1.3B parameters, and 2) Mistral-7B-Instruct-v0.2 (Jiang et al., 2023), a medium-sized model with 7B parameters. Both models were fine-tuned with LoRA (Hu et al., 2022), targeting the q proj, k proj, and v proj modules. The following LoRA parameters were used for fine-tuning: 1) rank of low-rank factorization (lora r) = 8, 2) scaling factor for the rank (lora alpha) = 32, and 3) lora dropout = 0.1. Additional fine-tuning parameters included: 1) learning rate = 1e-4, 2) weight decay = 0.05, and 3) batch size = 96 (for Mistral-7B-Instruct-v0.2) and 128 (for Phi-1.5).

3. Results and Discussion

Development of the shared vocabulary is a good example how partial "skills" of Large Language Models need to be mitigated to help them operate in the desired manner. LLMs have familiarity with SMILES notation and chemical structure concepts. They are neither consistent, nor generalizable, nor exhaustive.

To further assess the performance of LLMs in the Final referential chemical game, we used two language models: Phi-1.5 and Mistral-7B-Instruct-v0.2. For each LLM, we considered the base model with zero-shot and two-shot prompting techniques, as well as a fine-tuned model. Table 1 presents

LLM	Train			Test		
	Sender (Exact Match)	Sender (Partial Match)	Receiver (Exact Match)	Sender (Exact Match)	Sender (Partial Match)	Receiver (Exact Match)
Phi-1.5 zero-shot	0.0%	0.0%	0.0%	0.0%	0.0%	0.0%
Phi-1.5 2-shot	1.8%	3.5%	0.0%	0.0%	0.0%	0.0%
Phi-1.5 Fine-tuned	33.0%	71.3%	50.4%	0.0%	36.0%	32.1%
Mistral zero-shot	2.7%	36.3%	13.3%	3.4%	51.7%	6.5%
Mistral 2-shot	14.7%	66.5%	50.40%	12.1%	81.8%	15.8%
Mistral Fine-tuned	96.9%	99.6%	100.0%	72.2%	91.7%	68.6%

Table 1. Accuracy scores assessing Sender’s ability to construct Utterance from SMILES (World Object) and Receiver’s ability to reconstruct SMILES notation from Utterance. Exact Match measures if the Sender/Receiver’s output fully matched expected output. Partial Match for Sender counts instances where at least one part of Utterance is correct. Fine-tuned LLMs Phi-1.5 and Mistral-7B-Instruct-v0.2 shows significant improvement over base model with zero-shot and two-shot prompts

the results from various models for the referential game. We measured the accuracy of the Sender generating Utterance and the Receiver reconstructing SMILES separately. In the Train and Test games, the fine-tuned Mistral model significantly outperformed other models in Sender and Receiver accuracy with 72.2% and 68.6% respectively for test split.

The zero-shot and two-shot accuracy results for Phi-1.5 and Mistral models demonstrate the base models’ inability to parse and reason with SMILES notation of molecules. Mistral was able to understand SMILES better than the smaller Phi-1.5, as shown in the two-shot results. Fine-tuning with data created from the Molecular combinatorial library improved the capability of these models to understand, parse, and compose SMILES notation. Even after fine-tuning, Phi-1.5 was still unable to generate Utterance from SMILES, as indicated by the 0% Exact Match accuracy and only 36% Partial Match accuracy. However, Mistral handled SMILES notation much better after fine-tuning, with 72.2% and 91.7% accuracy in Exact Match and Partial Match, respectively.

We evaluate compositionality of the communication as topographic similarity (Brighton & Kirby, 2006; Lazaridou & Baroni, 2020; Korbak et al., 2019) - Spearman correlation of in-world distances between the objects (SMILES strings representing molecules) and their semantic distances. Semantic distances are evaluated as Cosine distances between embedding vectors of the names produced by the Sender. In-world distances are evaluated as Levenshtein editing distances between SMILES strings and Dice distances between Morgan fingerprints (Rogers & Hahn, 2010) of SMILES strings. Embeddings are obtained using all-MiniLM-L6-v2 sentence-transformer model (Reimers & Gurevych, 2019). With the base Mistral model (Mistral zero-shot), topographic similarity $\rho_{Levenshtein}$ is 0.07 and ρ_{Dice} is 0.09. Performance improvement of the fine-tuned model (Mistral Fine-tuned) is accompanied by appreciable increase of topographic similarity: $\rho_{Levenshtein}$ is 0.65 and

ρ_{Dice} is 0.82.

4. Conclusion

To our knowledge, this is the first attempt to leverage complex referential game setting to improve compositionality of communication between general-purpose LLMs.

It is tempting to consider RL-like setting of the referential game involving LLMs, where instead of fine-tuning (either in RL manner or supervised learning manner) the desired behavior is reinforced via prompting. Success of this approach appears to be highly sensitive to the nature of the LLM, just like with other prompt-driven reasoning strategies.

We would like to draw a deeper parallels with the field of emergent communication in MARL and notice that contemporary studies typically involve complex agent architectures with separate modules responsible for perception and communication. It seems that the demand for seamless communication with human agents calls for adoption of LLMs as enablers of shared grounding. Compositionality and reasoning, however, might be better delegated to the higher-level agents interacting with LLMs. In this case, the focus of communication games shifts from the emergence of language to the emergence of reasoning as a response to the complexity of the environment and interactions between agents.

References

- AI@Meta. Llama 3 model card. 2024. URL https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md.
- Barrett, J. A. and Skyrms, B. Self-assembling games. *The British Journal for the Philosophy of Science*, 68(2):329–353, 2017. doi: 10.1093/bjps/axv043. URL <https://doi.org/10.1093/bjps/axv043>.

- Brighton, H. and Kirby, S. Understanding Linguistic Evolution by Visualizing the Emergence of Topographic Mappings. *Artificial Life*, 12(2):229–242, 04 2006. ISSN 1064-5462. doi: 10.1162/artl.2006.12.2.229. URL <https://doi.org/10.1162/artl.2006.12.2.229>.
- Chaabouni, R., Kharitonov, E., Bouchacourt, D., Dupoux, E., and Baroni, M. Compositionality and generalization in emergent languages. In Jurafsky, D., Chai, J., Schluter, N., and Tetreault, J. (eds.), *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 4427–4442, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.407. URL <https://aclanthology.org/2020.acl-main.407>.
- Chen, J., Pan, X., Yu, D., Song, K., Wang, X., Yu, D., and Chen, J. Skills-in-context prompting: Unlocking compositionality in large language models. *ArXiv*, abs/2308.00304, 2023. URL <https://api.semanticscholar.org/CorpusID:260351132>.
- Fang, Y., Liang, X., Zhang, N., Liu, K., Huang, R., Chen, Z., Fan, X., and Chen, H. Mol-instructions: A large-scale biomolecular instruction dataset for large language models. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=Flsdsb619n>.
- Guo, T., Guo, K., Nan, B., Liang, Z., Guo, Z., Chawla, N. V., Wiest, O., and Zhang, X. What can large language models do in chemistry? a comprehensive benchmark on eight tasks. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023. URL <https://openreview.net/forum?id=lngbr3SZHW>.
- Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., and Chen, W. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=nZeVKeeFYf9>.
- Jiang, A. Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D. S., de Las Casas, D., Bressand, F., Lengyel, G., Lample, G., Saulnier, L., Lavaud, L. R., Lachaux, M., Stock, P., Scao, T. L., Lavril, T., Wang, T., Lacroix, T., and Sayed, W. E. Mistral 7b. *CoRR*, abs/2310.06825, 2023. doi: 10.48550/ARXIV.2310.06825. URL <https://doi.org/10.48550/arXiv.2310.06825>.
- Korbak, T., Zubek, J., Kucinski, L., Milos, P., and Raczaszek-Leonardi, J. Developmentally motivated emergence of compositional communication via template transfer. *ArXiv*, abs/1910.06079, 2019. URL <https://api.semanticscholar.org/CorpusID:204509150>.
- Kottur, S., Moura, J. M. F., Lee, S., and Batra, D. Natural language does not emerge ‘naturally’ in multi-agent dialog. In *Conference on Empirical Methods in Natural Language Processing*, 2017. URL <https://api.semanticscholar.org/CorpusID:6683636>.
- Lazaridou, A. and Baroni, M. Emergent multi-agent communication in the deep learning era. *ArXiv*, abs/2006.02419, 2020. URL <https://api.semanticscholar.org/CorpusID:219260403>.
- Lewis, D. K. *Convention: A Philosophical Study*. Wiley-Blackwell, Cambridge, MA, USA, 1969.
- Li, Y., Bubeck, S., Eldan, R., Del Giorno, A., Gunasekar, S., and Lee, Y. T. Textbooks are all you need ii: **phi-1.5** technical report. *arXiv preprint arXiv:2309.05463*, 2023.
- Reimers, N. and Gurevych, I. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 11 2019. URL <https://arxiv.org/abs/1908.10084>.
- Ristoski, P., Zubarev, D. Y., Gentile, A. L., Park, N., Sanders, D., Gruhl, D., Kato, L., and Welch, S. Expert-in-the-loop ai for polymer discovery. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management, CIKM '20*, pp. 2701–2708, New York, NY, USA, 2020. Association for Computing Machinery. ISBN 9781450368599. doi: 10.1145/3340531.3416020. URL <https://doi.org/10.1145/3340531.3416020>.
- Rogers, D. and Hahn, M. Extended-connectivity fingerprints. *Journal of Chemical Information and Modeling*, 50(5):742–754, 2010. doi: 10.1021/ci100050t. PMID: 20426451.
- SHAO, N., Cai, Z., xu, H., Liao, C., Zheng, Y., and Yang, Z. Compositional task representations for large language models. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=6axIMJA7ME3>.