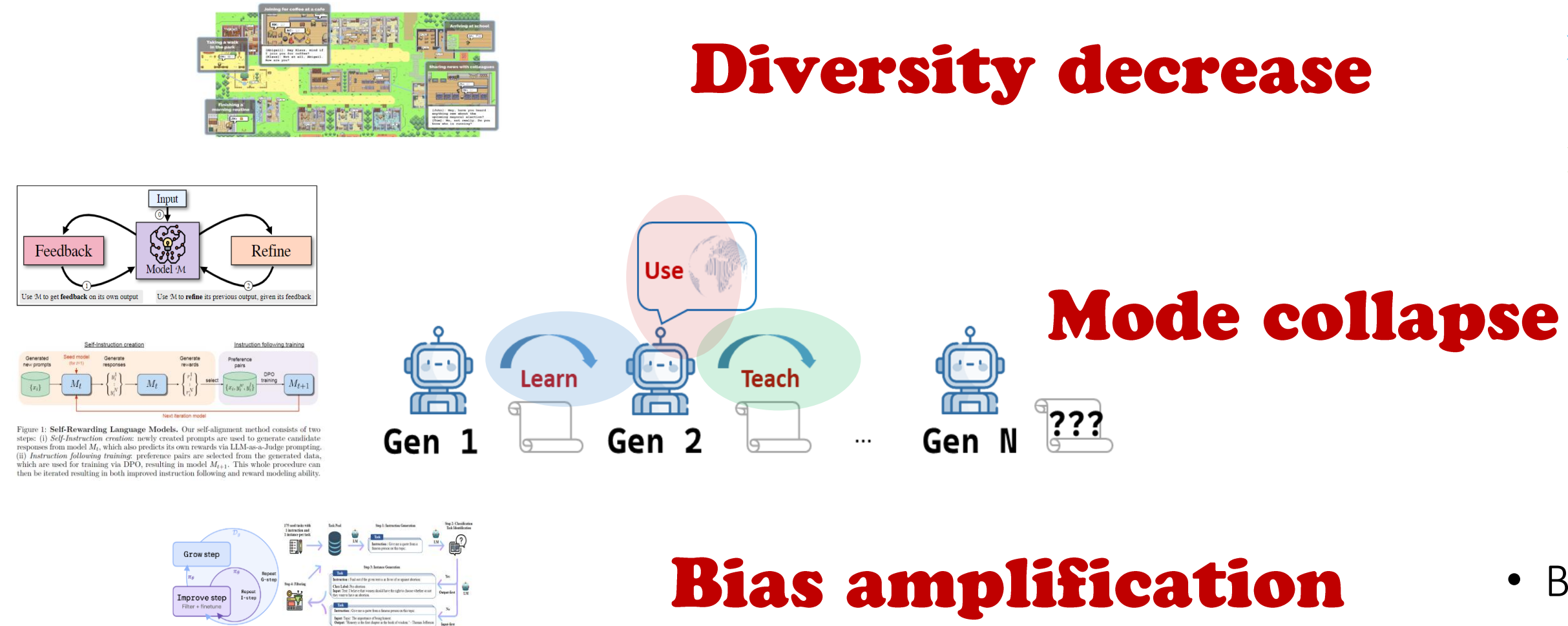


1. Motivation – what if self-improve too much?

- Self-interaction among LLM agents gains popularity, but **RISK?**

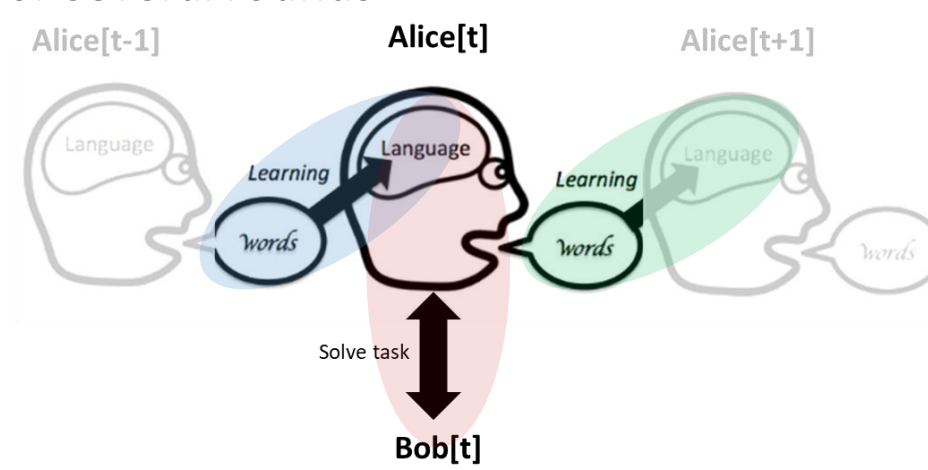


[1] Gulcehre, Caglar, et al. "Reinforced self-training (ReST) for language modeling." arXiv 2023.
 [2] Yuan, Weizhe, et al. "Self-rewarding language models." arXiv preprint arXiv 2024.
 [3] Madaan, Aman, et al. "Self-refine: Iterative refinement with self-feedback." NeurIPS 2023
 [4] Wang, Yizhong, et al. "Self-Instruct: Aligning Language Models with Self-Generated Instructions." ACL 2023
 [5] Gou, Zhibin, et al. "CRITIC: Large Language Models Can Self-Correct with Tool-Interactive Critiquing." ICLR 2024

2. Similarity to Human Language's Evolution

- Although proposed by different reasons, they are similar in:

- Imitation:** Another agent learn from the message generated by previous agent
- Interaction:** LLM interact with other or environment to refine the knowledge
- Transmission:** LLM generate message based on given prompts
- Repeat** the process for several rounds

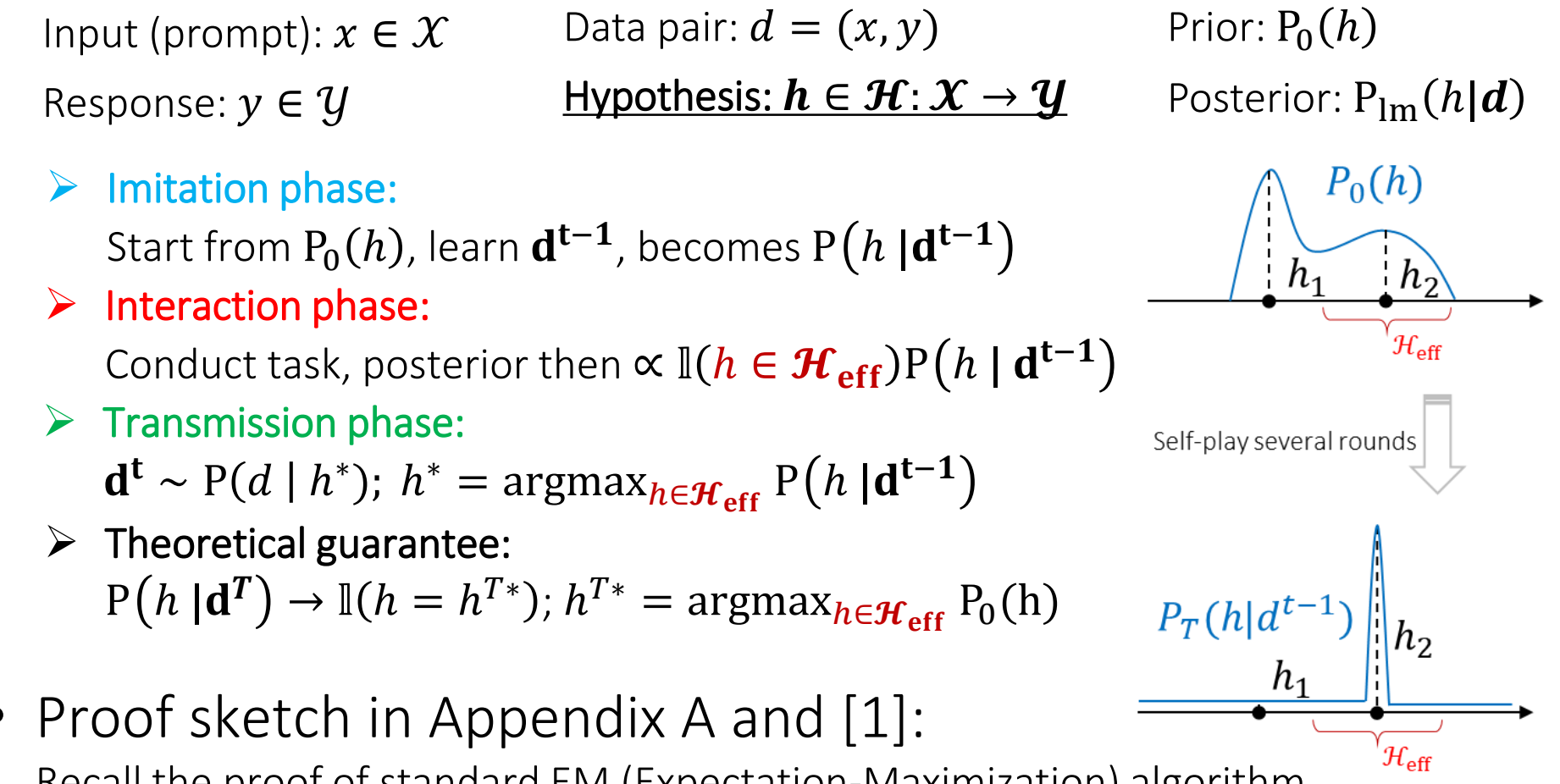


- But, keeps self-boosting introduce **RISKS**

Although reported in many related works sporadically, no **unified framework** to analyze the **asymptotic behavior**.

3. Bayesian Iterated Learning

- Bayesian-iterated learning framework:



- Input (prompt): $x \in \mathcal{X}$ Data pair: $d = (x, y)$
- Response: $y \in \mathcal{Y}$ Hypothesis: $h \in \mathcal{H}: \mathcal{X} \rightarrow \mathcal{Y}$
- Imitation phase:** Start from $P_0(h)$, learn \mathbf{d}^{t-1} , becomes $P(h | \mathbf{d}^{t-1})$
- Interaction phase:** Conduct task, posterior then $\propto \mathbb{I}(h \in \mathcal{H}_{\text{eff}})P(h | \mathbf{d}^{t-1})$
- Transmission phase:** $\mathbf{d}^t \sim P(d | h^*); h^* = \text{argmax}_{h \in \mathcal{H}_{\text{eff}}} P(h | \mathbf{d}^{t-1})$
- Theoretical guarantee:** $P(h | \mathbf{d}^T) \rightarrow \mathbb{I}(h = h^{T*}); h^{T*} = \text{argmax}_{h \in \mathcal{H}_{\text{eff}}} P_0(h)$

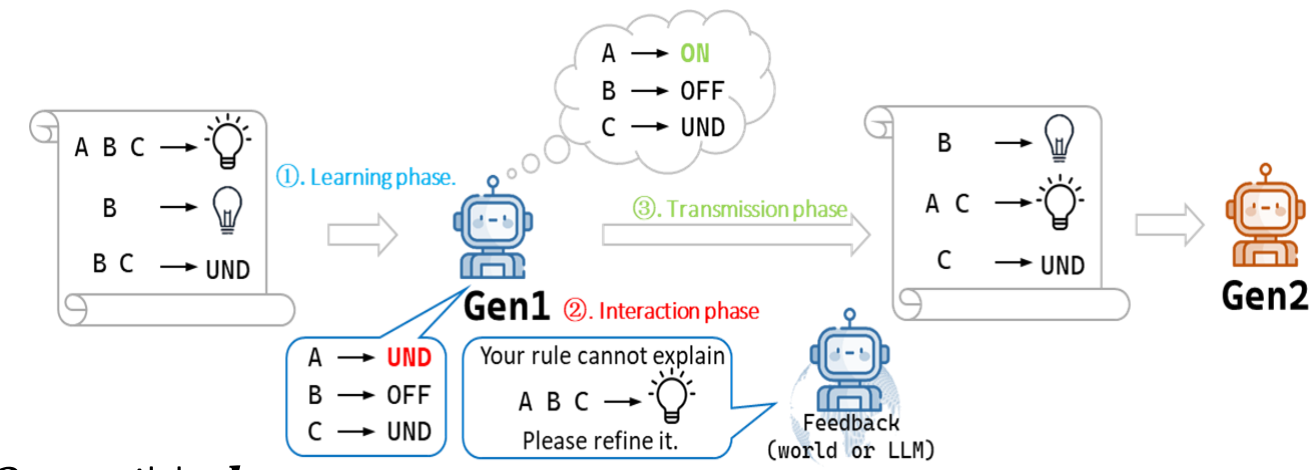
- Proof sketch in Appendix A and [1]: Recall the proof of standard EM (Expectation-Maximization) algorithm, replace (θ, z) to (h, \mathbf{d}) , marginalize the input variable x . Done!

Key assumption to LLM: ICL is implicit Bayesian Inference [2]

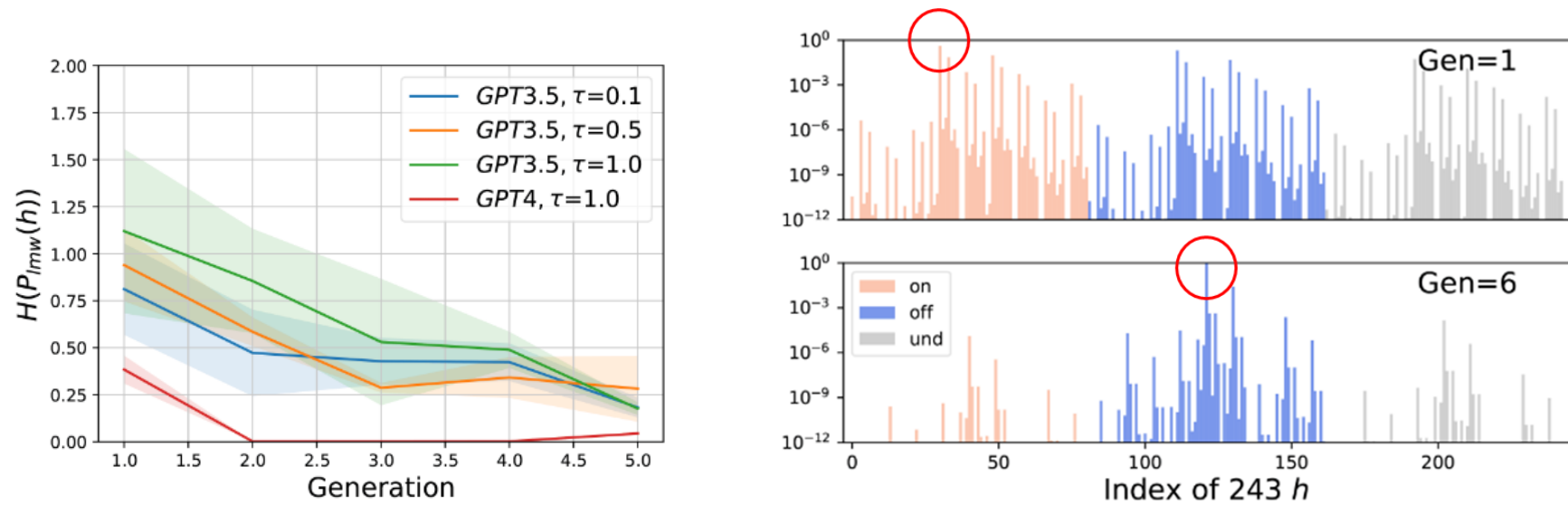
[1] Griffiths, Thomas L et al. "Using category structures to test iterated learning as a method for identifying inductive biases." Cognitive Science 2008.
 [2] Xie, Sang Michael, et al. "An Explanation of In-context Learning as Implicit Bayesian Inference." ICLR-2022

4. LLM Verification – Explicit h

- To verify the **subtle trends** predicted by the theory, start from Abstract Causal REasoning **ACRE**, used in [1]



- Consider 5 objects, then $3^5 = 243$ possible h

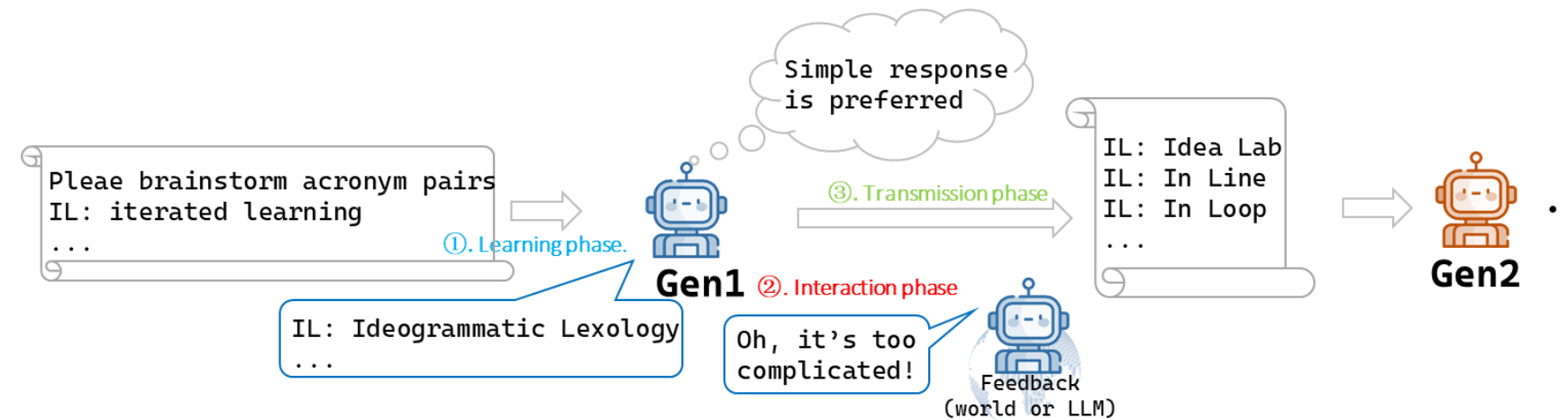


Verify convergence: $P(h | \mathbf{d}^t) \rightarrow \mathbb{I}(h = h^{T*})$ Verify solution: $h^{T*} = \text{argmax}_{h \in \mathcal{H}_{\text{eff}}} P_0(h)$

[1] Qiu, Linlu, et al. "Phenomenal Yet Puzzling: Testing Inductive Reasoning Capabilities of Language Models with Hypothesis Refinement." ICLR-2024

5. LLM Verification – Implicit h

- Consider a more practical self-data augmentation problem, where h is **implicit**, e.g., $h = \{\text{Long response, Short response}\}; h = \{\text{Use easy words, Use hard words}\}$
- A simple "creative writing"-style game, brainstorming the given acronym



- LLM naturally bias towards common & short words. Manipulate it using different \mathcal{H}_{eff}

Table 2: Results when adding different \mathcal{H}_{eff} . We color the **highest** and **lowest** numbers in each column. N_e is the number of easy examples in \mathbf{d}^0 . Results under different settings are in Table 4 and 5.

N_e	Ratio-easy				Avg-rank				Avg-length					
	2	4	6	8	2	4	6	8	2	4	6	8	10	
Random	0.91±0.01	0.60±0.08	0.96±0.00	0.87±0.03	0.82±0.06	13519	27269	7487	10425	15871	5.42±1.04	4.82±0.33	5.60±1.55	5.01±1.50
Imitation-only	0.43±0.20	0.93±0.01	0.92±0.00	0.97±0.00	0.96±0.00	35235	7497	9081	5549	8075	4.45±0.86	4.38±1.40	4.17±0.13	4.18±0.65
Hard	0.21±0.19	0.25±0.43	0.45±0.43	0.33±0.16	0.50±0.23	49869	46436	37288	41255	31903	4.62±1.54	5.78±1.39	4.67±0.40	4.38±0.60
Easy	0.76±0.17	1.00±0.00	0.98±0.00	1.00±0.00	0.99±0.00	15910	3156	2383	2924	2650	3.92±0.33	5.26±0.06	4.71±0.06	4.24±0.08
Easy-long	0.98±0.00	0.97±0.00	0.98±0.00	0.98±0.00	1.00±0.00	7063	9413	8649	6898	7404	5.20±0.41	5.88±0.52	6.38±1.10	6.97±1.57
Easy-short	1.00±0.00	1.00±0.00	0.97±0.00	1.00±0.00	0.98±0.00	5671	4223	5733	4502	5251	3.97±0.50	4.01±1.03	4.37±0.50	3.95±0.03

6. Take-away Message

- Applying Bayesian-IL to LLM's evolution:
 - Bias in $P_0(h)$ is guaranteed to be **amplified** if self-boosting **too much**
 - Bias can be **beneficial or harmful**, h can be explicit or implicit
 - Figure out the bias, understand it, and then design corresponding \mathcal{H}_{eff}
- Iterated learning and $P_0(h)$ in other fields:
 - CogSci: human prefer compositionality \rightarrow compositional language is achieved after IL
 - EmCom: simple NN prefer compositionality \rightarrow compositional mapping is achieved after IL
 - Representation Learning: complex NN prefer systematicness \rightarrow systematical generalization
 - VLM: language prefer compositionality \rightarrow vision modal also becomes compositional after IL
- In-weights updates (e.g., DPO) amplify the bias in $P_0(h)$ more
- Analysis of the "squeezing effect" caused by negative gradient part in DPO

[1] Kirby, Simon, et al. "Cumulative cultural evolution in the laboratory: An experimental approach to the origins of structure in human language." PNAS 2008
 [2] Ren, Yi, et al. "Compositional languages emerge in a neural iterated learning model." ICLR 2020
 [3] Ren, Yi, et al. "Improving compositional generalization using iterated learning and simplicial embeddings." NeurIPS 2023
 [4] Zheng, Chenhao, et al. "Iterated learning improves compositionality in large vision-language models." CVPR 2024
 [5] Ren, Yi, et al. "Learning Dynamics of LLM Finetuning", Submitted to ICLR 2025