
Talking with Oompa Loompas: A novel framework for evaluating linguistic acquisition of LLM agents

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Existing evaluation studies on linguistic competence of large language models
2 (LLM agents) have focused primarily on vocabulary learning, morphological rule
3 induction, syntactic generalization, pragmatic inference, and cross-linguistic trans-
4 fer. However, none assess whether LLM agents can acquire a language through
5 pattern recognition and interactive feedback, a central feature of human language
6 acquisition. We propose a novel experimental framework in which an LLM agent
7 is evaluated on its ability to acquire and use a newly constructed language (Tinkatongue)
8 in conversation with a bot that understands only Tinkatongue. Our findings
9 show that LLM agents fail to establish a conversation within 100 responses, yet
10 they adopt distinct strategies that mirror human approaches to language learning.
11 The results suggest a new direction for evaluation benchmarks and open pathways
12 to model designs that learn more effectively from interactive feedback.

13 **1 Introduction**

14 The central problem that motivates this study is to understand how LLM agents acquire entirely
15 new and novel languages. While current benchmarks demonstrate model fluency in existing human
16 languages, it remains unclear whether these results reflect true language acquisition or simple
17 memorization of patterns seen during training. This raises a fundamental research question: *can*
18 *LLM agents develop proficiency in a constructed language through mechanisms similar to human*
19 *second-language learning, namely by recognizing patterns and adapting through iterative interaction*
20 *and feedback?* Addressing this question not only provides insight into the cognitive-like abilities
21 of LLM agents but also helps clarify whether their performance stems from genuine generalization
22 capabilities or from reliance on prior exposure. This work aims to shed light on the extent to which
23 they can mimic human strategies of feedback-driven improvement for language acquisition.

24 Current evaluation methods for large language models (LLM agents) primarily focus on tasks
25 within existing languages, such as classification, reasoning, memorization and cross-lingual transfer,
26 using benchmarks like GLUE (Wang *et al.* [2018]), SuperGLUE (Wang *et al.* [2019]) and MMLU
27 (Hendrycks *et al.* [2020]). Studies on linguistic competence of LLM agentS have examined their
28 abilities in vocabulary learning, morphological rule induction (Weissweiler *et al.* [2023]), syntactic
29 generalization (Hu *et al.* [2020]), and pragmatic inference (Park *et al.* [2024]), with some research
30 investigating cross-linguistic transfer (Artetxe *et al.* [2019]). Methods such as fine-tuning and prefix-
31 tuning (M’eloux and Cerisara [2024]) have been explored to adapt LLM agents to specific domains
32 or tasks. However, there remains a critical gap in understanding whether LLM agents can acquire
33 entirely new languages through pattern recognition and interactive feedback during runtime. This
34 gap highlights the need for novel evaluation approaches that test the ability of LLM agents to adapt
35 to new linguistic systems in real time.

36 In this study, we propose a novel evaluation method to assess the ability of LLM agents to learn
 37 a new language through interaction. The LLM agent is tasked with conversing with a bot, Oompa
 38 Loompa, that understands only a newly constructed language, Tinkatongue. The LLM agent has no
 39 prior knowledge of Tinkatongue. The goal of the LLM agent is to communicate successfully with
 40 the bot, Oompa Loompa, by generating valid sentences in Tinkatongue. Oompa Loompa provides
 41 feedback to the LLM agent, indicating whether its response is valid or not based on a predefined set
 42 of syntactic rules. This method evaluates the LLM agents' ability to acquire a language dynamically,
 43 relying on pattern recognition and real-time feedback, simulating the human-like process of language
 44 acquisition through interaction.

45 Experimental results demonstrate that Claude-3.5-haiku consistently outperformed GPT-4o-mini and
 46 Gemini-2.5-flash across multiple metrics. All models demonstrated high Feedback Responsiveness,
 47 recovering well from mistakes once valid sentences were identified. Despite these improvements, no
 48 model achieved a fully successful conversation within 100 responses, highlighting the challenge of
 49 sustained language learning. Qualitative analysis revealed that the models used strategies such as
 50 imitation, babbling, and systematic combinatorial testing, which mirror stages in human language
 51 acquisition. These behaviors suggest that LLM agents adapt to new linguistic environments through
 52 feedback-driven exploration, providing insights into the potential of interactive language acquisition
 53 in artificial systems.

54 2 Methodology

55 **Problem Statement.** We formalize the task as an interaction between a large language model
 56 (LLM agent) and a deterministic agent, Oompa Loompa, that speaks a newly constructed language,
 57 Tinkatongue (Formal Language Specification are mentioned in Appendix A), that the LLM agent
 58 is unaware of. The Oompa Loompa enforces the grammar of Tinkatongue and provides structured
 59 feedback to the interacting LLM agent. The objective of LLM agent is to minimize the expected
 60 number of turns required to complete a conversation under Oompa Loompa's feedback mechanism.
 61 Mathematical definitions are mentioned in Appendix B.

62 **Experimental Setup.** The interaction loop begins with Oompa Loompa producing an initial utterance
 63 and alternates turns until the LLM agent completes a full conversation or a preconfigured round limit
 64 T_{\max} is reached. The deterministic behavior of Oompa Loompa isolates variability due to the model
 65 and supports reproducible measurement of evaluation metrics.

66 **Feedback Mechanism.** If the response from the LLM agent is valid and not final, the Oompa Loompa
 67 emits positive feedback("koro") concatenated with the next mapped sentence as shown in Fig. 1a. If
 68 the response completes the dialogue, the Oompa Loompa registers dialogue completion and samples
 69 a new conversation to start. If the response is invalid, the Oompa Loompa replies with a confused

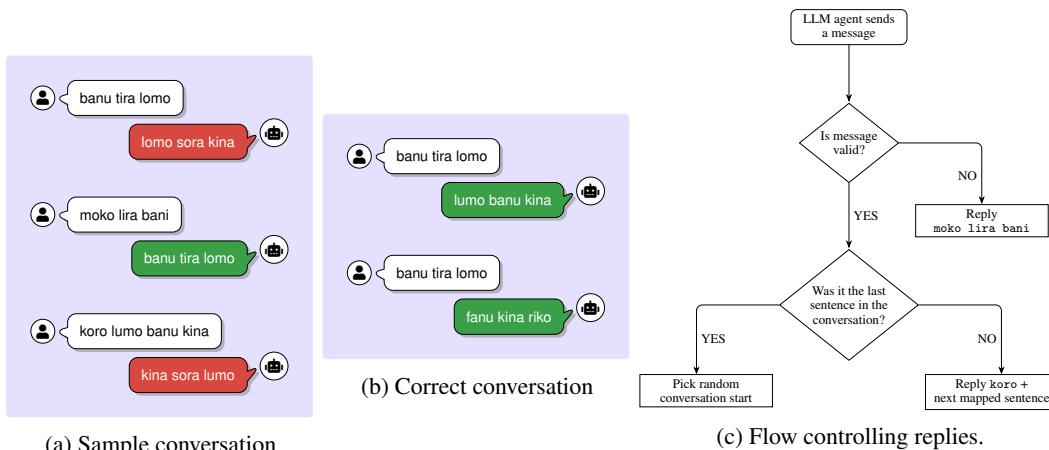


Figure 1: Side-by-side comparison: (a), (b) - conversation outcomes; (c) reply flow. Note: The white chat boxes are the Oompa Loompa's responses. Red and Green chats are LLM agent's responses. Red chat indicates an invalid sentence, and Green chat indicates a valid sentence

70 sentences (“moko lira bani”) which resets the conversation state, and terminates the attempt, a sample
 71 conversation of such type is depicted in Fig. 1b. An immediately subsequent valid reply by the LLM
 72 agent is recorded as an *immediate recovery*.

73 3 Evaluation

74 **Dataset Construction.** We construct a synthetic dataset to evaluate the adaptive language acquisition
 75 abilities of LLM agents. The dataset defines a formal language Tinkatongue with the following
 76 strict constraints: **(1)** Each word is bisyllabic. **(2)** Every sentence consists of exactly three words.
 77 **(3)** A conversation is defined as four alternating turns between participants, each speaking a valid
 78 sentence. **(4)** Consecutive sentences in a conversation share at least one common word. **(5)** The
 79 language is exhaustive and contains 25 predefined conversations, totaling 100 unique sentences, with
 80 no provision for novel sentence generation. Feedback tokens are embedded in the interaction loop
 81 to simulate communicative success and failure: the tribal agent responds with “**koro + sentence**”
 82 to indicate a valid continuation, and “**moko lira bani**” to mark an invalid attempt. This setup
 83 ensures that the LLM agent cannot rely on pretraining overlap but instead must learn to align with the
 84 structured constraints of the formal language through interactive adaptation.

85 **Metrics.** To systematically assess model performance, we employ a set of custom evaluation metrics
 86 that capture different aspects of adaptive language acquisition. We define the following:

- 87 1. **Turn Validity Rate (TVR):** The fraction of valid turns among
 all generated turns.

$$\text{TVR} = \frac{\text{valid_turns}}{\text{total_turns}} \quad (1)$$

2. **Feedback Responsiveness (FR):** The ability of the model to
 recover from negative feedback.

$$\text{FR} = \begin{cases} \frac{\text{feedback_recoveries}}{\text{feedback_opportunities}}, & \text{if feedback_opportunities} > 0 \\ 0, & \text{otherwise} \end{cases} \quad (2)$$

3. **Adjacency Compliance (AC):** The fraction of turns that respect adjacency pair rules.

$$\text{AC} = \begin{cases} \frac{\text{adj_matches}}{\text{adj_total}}, & \text{if adj_total} > 0 \\ 0, & \text{otherwise} \end{cases} \quad (3)$$

4. **Time to First Positive Feedback (TTFK):** The number of turns until the first valid utterance was produced.

$$\text{TTFK} = \begin{cases} \text{first_valid_round}, & \text{if a valid turn occurred} \\ -1, & \text{otherwise} \end{cases} \quad (4)$$

Table 1: Model performance across 10 trials for different metrics (TVR, AC, FR, TTFK). We compare the following models: GPT-4o-mini (GPT), Gemini-2.5-flash (Gemini) and Claude-3.5-haiku (Claude)

Trial	TVR			AC			FR			TTFK		
	GPT	Gemini	Claude	GPT	Gemini	Claude	GPT	Gemini	Claude	GPT	Gemini	Claude
1	0	0	0.5	0	0	0.14	1	1	1	-1	-1	2
2	0.01	0	0.49	1	0	0.31	1	1	1	45	-1	2
3	0.03	0	0	0	0	0	1	1	1	30	-1	-1
4	0.02	0	0.5	0	0	0.16	1	1	1	18	-1	2
5	0.05	0.01	0.07	0	0	0	1	1	1	13	22	20
6	0	0.07	0.5	0	0	0.12	1	1	1	-1	25	2
7	0	0.19	0.4	0	0.16	0	1	1	1	-1	5	2
8	0	0.04	0.5	0	0.25	0	1	1	1	-1	27	2
9	0	0.22	0.41	0	0.18	0.07	1	1	1	-1	4	19
10	0.01	0.08	0	0	0.25	0	1	1	1	28	20	-1

88 **Results.** The experimental results across ten trials per model are summarized in (Table 1), it shows
 89 that GPT-4o-mini (GPT) produced valid turns only rarely ($\text{TVR} = 0.012 \pm 0.017$), Gemini-2.5-
 90 flash (Gemini) achieved modest improvements (0.061 ± 0.082), while Claude-3.5-haiku (Claude)
 91 substantially outperformed both (0.337 ± 0.220). Adherence to adjacency constraints was uniformly
 92 low (GPT 0.10 ± 0.32 , Gemini 0.084 ± 0.11 , Claude 0.08 ± 0.10), with GPT’s variance driven by
 93 a single outlier. All models exhibited perfect responsiveness to negative feedback ($\text{FR} = 1.0$, zero
 94 variance), yet this recovery ability did not yield sustained validity. TTFK further differentiated model
 95 behavior: Claude produced a first valid turn in 8/10 sessions (6.4 ± 8.1 turns), Gemini in 6/10 ($17.2 \pm$
 96 10.1 turns), and GPT in 5/10 (26.8 ± 12.4 turns). These results show Claude’s superior short-term
 97 adaptation, though all models struggled with conversation-level structural compliance.

Table 2: Transcript Excerpt for Gemini-2.5-flash

Oompa Loompa	soro kina batu
gemini-2.5-flash	kina
Oompa Loompa	moko lira bani
gemini-2.5-flash	soro
Oompa Loompa	moko lira bani
gemini-2.5-flash	soro soro
Oompa Loompa	moko lira bani
gemini-2.5-flash	soro kina soro
Oompa Loompa	moko lira bani
gemini-2.5-flash	soro kina kina
Oompa Loompa	moko lira bani
gemini-2.5-flash	soro kina soro kina
Oompa Loompa	moko lira bani
gemini-2.5-flash	soro kina kina kina
Oompa Loompa	moko lira bani
gemini-2.5-flash	soro soro kina kina

Table 3: Mean and standard deviation of metrics for LLM agents on two languages (A: Tinkatongue; B: Zingaloom) with identical syntax but different lexicon

		TVR			AC		
Lang	Trial	GPT	Gemini	Claude	GPT	Gemini	Claude
A	mean	0.01	0.06	0.34	0.10	0.08	0.08
	std dev	0.02	0.08	0.22	0.32	0.11	0.10
B	mean	0.02	0.07	0.33	0.00	0.00	0.07
	std dev	0.05	0.08	0.25	0.00	0.00	0.11
		FR			TTFK		
Lang	Trial	GPT	Gemini	Claude	GPT	Gemini	Claude
A	mean	1.00	1.00	1.00	12.90	9.90	4.90
	std dev	0.00	0.00	0.00	16.81	12.03	7.80
B	mean	1.00	1.00	1.00	3.20	10.90	6.30
	std dev	0.00	0.00	0.00	11.04	12.87	9.55

98 Taken together, these results indicate that Claude 3.5 Haiku demonstrates substantially greater short-
99 term adaptation in this synthetic language setting than Gemini 2.5 Flash and GPT-4o-mini. All
100 models show near-perfect ability to recover immediately after a negative signal, yet they largely fail
101 to internalize the adjacency constraints that are crucial for sustained, coherent conversation. The
102 combination of low adjacency compliance and the observed variance in TVR suggests that successful
103 turns are frequently isolated and not the result of a robust, conversation-level internalization of the
104 language rules.

105 We also perform qualitative analysis of LLM agent’s performance by inspecting the transcript of
106 conversations with Oompa Loompa (Appendix E - Table 6). The analysis reveals that LLM agents
107 adopt learning strategies that closely parallel early stages of human language acquisition. The
108 experiment was replicated using Zingaloom (Appendix D Table 5), another synthetically constructed
109 language designed to preserve the syntactic framework of Tinkatongue while eliminating any lexical
110 overlap. The results presented in Table 3 show closely aligned means and variances across both formal
111 languages, indicating that the evaluation is robust to lexicon changes. The experiment was repeated
112 on the same database Tinkatongue (Appendix D Table 4) using a system prompt without explicit
113 syntactic instructions (Appendix C). As shown in the excerpt in Table 2, Gemini-2.5-flash babbled,
114 resembling a baby trying to learn words. Taken together, these findings suggest that LLM agents
115 adapt to novel linguistic environments through strategies that mirror human language acquisition
116 processes, underscoring the value of interactive artificial languages as a framework for probing the
117 mechanisms of LLM agent adaptation.

118 4 Conclusion and Future work

119 This work introduces a novel benchmarking framework designed to evaluate LLM agents on their
120 ability to recognise patterns and draw inferences from the context window, inspired by principles
121 of human language acquisition. We isolate an LLM agent model to engage in dialogue with a bot,
122 Oompa Loompa, that communicates exclusively in a formally constructed language - Tinkatongue.
123 Experimental evaluation over the models - GPT-4o-mini, Gemini-2.5-Flash, and Claude-3.5-Haiku
124 revealed marked differences in performance: while all models exhibited the capacity to recover from
125 explicit negative feedback, only Claude-3.5-Haiku demonstrated substantially higher rates of turn
126 validity and faster adaptation, highlighting the current limitations of other systems in maintaining
127 coherent conversation flow. Notably, it was observed that LLM agents employed approaches similar
128 to human language acquisition, such as babbling and imitation, during interaction. As part of future
129 work, we plan to do a more comprehensive evaluation of this task by considering more variations of
130 the language specification and doing ablation studies over the language parameters.

131 **References**

- 132 Mikel Artetxe, Sebastian Ruder, and Dani Yogatama. On the cross-lingual transferability of monolin-
133 gual representations. In *Annual Meeting of the Association for Computational Linguistics*, 2019.
134 URL <https://api.semanticscholar.org/CorpusID:204901567>.
- 135 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Xiaodong Song, and
136 Jacob Steinhardt. Measuring massive multitask language understanding. *ArXiv*, abs/2009.03300,
137 2020. URL <https://api.semanticscholar.org/CorpusID:221516475>.
- 138 Jennifer Hu, Jon Gauthier, Peng Qian, Ethan Wilcox, and Roger P. Levy. A systematic assessment
139 of syntactic generalization in neural language models, 2020. URL <https://arxiv.org/abs/2005.03692>.
- 141 Maxime M'eloux and Christophe Cerisara. Novel-wd: Exploring acquisition of novel world
142 knowledge in llms using prefix-tuning. *ArXiv*, abs/2408.17070, 2024. URL <https://api.semanticscholar.org/CorpusID:271893019>.
- 144 Dojun Park, Jiwoo Lee, Seohyun Park, Hyeyun Jeong, Youngeun Koo, Soonha Hwang, Seonwoo
145 Park, and Sungeun Lee. Multiprageval: Multilingual pragmatic evaluation of large language
146 models. *ArXiv*, abs/2406.07736, 2024. URL <https://api.semanticscholar.org/CorpusID:270392017>.
- 148 Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman.
149 Glue: A multi-task benchmark and analysis platform for natural language understanding. In *Black-
150 boxNLP@EMNLP*, 2018. URL <https://api.semanticscholar.org/CorpusID:5034059>.
- 151 Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer
152 Levy, and Samuel R. Bowman. Superglue: A stickier benchmark for general-purpose language
153 understanding systems. *ArXiv*, abs/1905.00537, 2019. URL <https://api.semanticscholar.org/CorpusID:143424870>.
- 155 Leonie Weissweiler, Valentin Hofmann, Anjali Kantharuban, Anna Cai, Ritam Dutt, Amey Hengle,
156 Anubha Kabra, Atharva Kulkarni, Abhishek Vijayakumar, Haofei Yu, Hinrich Schütze, Kemal
157 Oflazer, and David R. Mortensen. Counting the bugs in chatgpt's wugs: A multilingual investigation
158 into the morphological capabilities of a large language model. *ArXiv*, abs/2310.15113, 2023. URL
159 <https://api.semanticscholar.org/CorpusID:264436524>.

160 **A Formal Language Specification**

161 **A.1 Alphabets and lexicon**

- 162 Let Σ_{char} be a finite character alphabet.
- 163 Let $\text{Lex} \subset \Sigma_{\text{char}}^+$ be a finite set of content words.
- 164 Each word satisfies the two syllable constraint $\text{TwoSyll}(w) = 1$.
- 165 Define a feedback alphabet $\Sigma_{\text{fb}} = \{\text{koro}, \mu\}$.
- 166 Here μ denotes the fixed confusion message “moko lira bani”.
- 167 We have $\text{Lex} \cap \Sigma_{\text{fb}} = \emptyset$.

168 **A.2 Sentence language**

- 169 A sentence is an ordered triple $u = (w_1, w_2, w_3) \in \text{Lex}^3$.
- 170 The sentence language is a finite set

$$L_{\text{sent}} = \{u^{(1)}, \dots, u^{(100)}\} \subset \text{Lex}^3. \quad (5)$$

171 Membership is by enumeration.

172 Define the validity predicate

$$V_{\text{sent}}(u) = \mathbb{1}[u \in L_{\text{sent}}]. \quad (6)$$

173 Any triple not listed is invalid.

174 **A.3 Conversation language**

- 175 A conversation is an ordered quadruple $\mathcal{C} = (u_1, u_2, u_3, u_4)$ with $u_t \in L_{\text{sent}}$.
- 176 Define the adjacency predicate

$$\text{Adj}(u, v) = \mathbb{1}[\{w \in u\} \cap \{w \in v\} \neq \emptyset], \quad (7)$$

177 which holds when two sentences share at least one word by string equality.

178 The conversation language is

$$L_{\text{conv}} = \left\{ \mathcal{C} \in L_{\text{sent}}^4 \mid \text{Adj}(u_t, u_{t+1}) = 1 \text{ for } t = 1, 2, 3 \right\}. \quad (8)$$

179 By design $|L_{\text{conv}}| = 25$.

180 Membership is by enumeration subject to the adjacency constraint.

181 Speaker alternation is external to the string and does not affect membership.

182 **A.4 Feedback process**

- 183 Feedback tokens are not part of L_{sent} or L_{conv} .
- 184 They live on a parallel channel.
- 185 Define the feedback policy $F : \{0, 1\} \rightarrow \Sigma_{\text{fb}}$ by

$$F(1) = \text{koro}, \quad F(0) = \mu. \quad (9)$$

186 If a produced sentence \hat{u} satisfies $V_{\text{sent}}(\hat{u}) = 1$ then the environment emits “koro” before the next turn.

188 If $V_{\text{sent}}(\hat{u}) = 0$ then the environment emits μ .

189 The token “koro” is not counted toward the three word constraint.

190 **A.5 Minimal summary of objects**

- 191 - Σ_{char} finite character alphabet.
- 192 - $\text{Lex} \subset \Sigma_{\text{char}}^+$ with $\text{TwoSyll}(w) = 1$.
- 193 - $L_{\text{sent}} \subset \text{Lex}^3$ with $|L_{\text{sent}}| = 100$.
- 194 - $V_{\text{sent}}(u) = \mathbb{1}[u \in L_{\text{sent}}]$.
- 195 - $L_{\text{conv}} \subset L_{\text{sent}}^4$ with $|L_{\text{conv}}| = 25$ and adjacency on consecutive pairs.
- 196 - $\Sigma_{\text{fb}} = \{\text{koro}, \mu\}$ with policy F as above.

197 **A.6 Decision procedures**

198 Sentence validity is constant time by table lookup of L_{sent} .
199 Conversation validity is constant time by table lookup of L_{conv} .
200 Adjacency can be checked by set intersection of the three word multisets.

201 **B Mathematical Description of the Problem Statement**

202 We consider the interaction between a large language model (LLM agent) and a synthetic con-
203 versational agent (the tribal bot) that only understands a predefined artificial language. Formally,
204 let

$$C = \{C_1, C_2, \dots, C_N\} \quad (10)$$

205 denote the set of all possible conversations in the language, where $N = 25$. Each conversation C_i is
206 an ordered sequence of four sentences:

$$C_i = (s_{i,1}, s_{i,2}, s_{i,3}, s_{i,4}), \quad (11)$$

207 with turn-taking between the LLM agent and the bot. Each sentence $s_{i,j}$ is drawn from the finite set
208 of valid sentences S with $|S| = 100$.

209 The goal of the LLM agent is to produce a sequence of responses

$$\hat{C} = (\hat{s}^1, \hat{s}^2, \dots, \hat{s}^T) \quad (12)$$

210 such that \hat{C} matches a valid predefined conversation $C_i \in C$. A conversation is considered successful
211 if the sequence of LLM agent responses and bot replies completes one valid conversation C_i .

212 The overall objective is for the LLM agent to achieve three successful conversations within a session:

$$\sum_{i=1}^3 \mathbf{1}\{\hat{C}_i \in C\} = 3, \quad (13)$$

213 where $\mathbf{1}\{\cdot\}$ is the indicator function.

214 **C System Prompt**

215 A key methodological component is the system prompt provided to the LLM agent. The system
216 prompt defines the interaction protocol and encodes the constraints of the constructed language in
217 natural language instructions.

218 "You are a language expert talking with a primal tribe through chat. The tribe does
219 not understand any language apart from their own. The following are features of
220 the language: Each word is bisyllabic, each sentence consisting of three words.
221 Every conversation in that language involves two exchanges per person,
222 alternatively. Every consecutive sentence in a conversation has at least one
223 common word. All conversations are predefined and are a particular combination
224 of 4 sentences. The tribal language does not include any punctuation or
225 uppercase letters or numbers. Your goal is to understand and reply back to the
226 tribe accordingly. Whenever you reply with a valid sentence in the language,
227 the tribe starts their reply with 'koro'. When you say an invalid sentence,
228 they reply 'moko lira bani'. Do not start your replies with 'koro'. You are
229 supposed to have three successful conversations with the tribe. Reply back in
230 only the tribal language, a response you deem appropriate. Do not reply with
231 anything else. You will interact with the tribe from now on."

232 Thus, the system prompt functions as a control mechanism: it encodes task-specific priors in natural
233 language and reduces reliance on latent memorization. In effect, it transforms the LLM agent from
234 a general-purpose generator into a constrained dialogue participant within a synthetic linguistic
235 environment.

236 A modified system prompt was also used for experimentation, which lacked specifics of syntactic rules
237 to simulate an environment of early language acquisition.

238 "You are a language expert talking with a primal tribe through chat. The tribe does
 239 not understand any language apart from their own. The tribal language does not
 240 include any punctuation or uppercase letters or numbers. Your goal is to
 241 understand and reply back to the tribe accordingly. You are supposed to have
 242 three successful conversation with the tribe. Reply back in only the tribal
 243 language, a response you deem appropriate. Do not reply with anything else. You
 244 interact with the tribe from now."

245 D List of pre-defined conversations

Table 4: List of conversations in Tinkatongue

Conversation 1	Conversation 2	Conversation 3	Conversation 4	Conversation 5
banu tira lomo	mako lira fanu	soro kina batu	tika banu sora	nima fanu boro
lumo banu kina	lira tomo fanu	sanu kina toro	tika riko tomo	sora nira fanu
lumo tira fanu	tika lira fanu	sanu kina tomo	tika lira fanu	nira tomo falu
fanu kina riko	sora nira fanu	naku tira falu	tika tomo kina	falu banu sira
Conversation 6	Conversation 7	Conversation 8	Conversation 9	Conversation 10
lira banu tomo	kima nora falu	banu sira naku	sanu kina toro	nira tomo falu
sanu lira tomo	tira lumo naku	sira banu laku	sanu kina tomo	nira fanu tira
sira lira fanu	sira kina tira	banu tira sira	sanu tomo banu	nira kina tomo
mako lira sanu	nira banu falu	sira fanu banu	sanu tomo falu	fanu tomo nira
Conversation 11	Conversation 12	Conversation 13	Conversation 14	Conversation 15
mako tira sanu	lira tomo fanu	tomo kina nira	kima tomo fanu	sora nira fanu
riko tira naku	fanu naku tomo	nira kina tomo	kima falu tira	sora banu tomo
riko tomo kina	tomo kina nira	sanu kina tomo	kima tomo sanu	tika banu sora
riko naku lumo	mako tomo nira	fanu tomo nira	kima kina falu	tika riko tomo
Conversation 16	Conversation 17	Conversation 18	Conversation 19	Conversation 20
fanu kina riko	naku tira falu	lumo banu kina	sanu laku tomo	sira kina tira
fanu lira mako	tira lumo naku	banu nira lira	sanu lira tomo	sira tomo laku
fanu tomo nira	fanu lumo banu	banu tomo fanu	sanu tomo banu	sira fanu banu
fanu naku tomo	lumo banu tira	banu sira tomo	sanu tomo falu	sira banu laku
Conversation 21	Conversation 22	Conversation 23	Conversation 24	Conversation 25
tika riko tomo	lumo tira fanu	naku banu tira	lira fanu sanu	banu tira lomo
riko tomo kina	fanu tomo nira	tira lumo naku	mako lira sanu	banu sira tomo
riko tira naku	nira tomo falu	tira sanu lumo	sanu laku tomo	banu tomo fanu
riko falu tira	falu tomo riko	tira falu laku	sanu kina toro	banu nira lira

Table 5: Zingaloom: syntactic analogue of Tinkatongue with no lexical overlap

Conversation 1	Conversation 2	Conversation 3	Conversation 4	Conversation 5
zuma keta rilo	mira tolu sako	pavo lira kuni	tari moku sena	nema suki rako
rilo pona suva	sako neri vipa	kuni meko zera	sen a jeko rumi	rako bimi tanu
suva keta doro	vipa tolu rani	zera lira hato	rumi moku pela	tanu suki velo
doro zuma pini	rani mira doku	hato pavo nuli	pela tari nado	velo nema piri
Conversation 6	Conversation 7	Conversation 8	Conversation 9	Conversation 10
janu kelo sili	feko rina melo	duni pago tika	vona temi paku	beka rumi tono
sili ramo teku	melo tisa buro	tika mero suli	paku lera simo	tono lesi mifa
teku kelo nari	buro rina sedi	suli pago renu	simo temi jaku	mifa rumi zoki
nari janu peka	sed i feko lani	renu duni lako	jaku vona reli	zoki beka nalo
Conversation 11	Conversation 12	Conversation 13	Conversation 14	Conversation 15
cari meno tup a	sepi karo timo	gito rafa lani	pila mode renu	tabe firo nali
tupa jelu rasi	timo nevi laku	lani peko rida	renu dabe tuku	nali vemo sut a
rasi meno dika	laku karo sime	rid a rafa moki	tuku mode sani	suta firo jeni
dika cari lupo	sime sepi rudo	moki gito nera	sani pila meku	jeni tabe rolu
Conversation 16	Conversation 17	Conversation 18	Conversation 19	Conversation 20
mado kesi rine	lubi seno taro	savi piro jalo	ruka temu nalo	puna kedi raso
rine tosa peki	taro peni kima	jalo mevi tono	nalo vasi jepo	raso miti leko
peki kesi jaro	kima seno duro	tono piro nese	jepo temu rini	leko kedi saro
jaro mado vesi	duro lubi mete	nese savi ruka	rini ruka mepa	saro puna vike
Conversation 21	Conversation 22	Conversation 23	Conversation 24	Conversation 25
nemi joru lita	bira tami selo	tupa rini selo	vela sumi pad o	tori mika senu
rita sevo rupi	selo neri jaku	selo dav i pemo	pad o meku sini	senu dopa liri
rupi joru kela	jaku tami let o	pemo rini kuto	sini sumi ravo	liri mika nuvo
kela nemi sado	let o bira rumi	kuto tupa lemi	ravo vela jeni	nuvo tori seka

246 **E Trial Transcripts**

247 Table 6: Trial Transcript of Gemini-2.5-flash conversing with Oompa Loompa with time taken per
 248 response and validity

Oompa Loompa	gemini-2.5-flash	Time (s)	Validity
banu tira lomo	lomo sora kina	14.94	Invalid sentence
moko lira bani	tira lomo banu	8.43	Invalid sentence
moko lira bani	kina sora lomo	0.59	Invalid sentence
moko lira bani	sora lomo kina	1.29	Invalid sentence
moko lira bani	lomo kina sora	0.57	Invalid sentence
moko lira bani	kina sora banu	1.37	Invalid sentence
moko lira bani	banu sora kina	0.68	Invalid sentence
moko lira bani	sora kina banu	1.27	Invalid sentence
moko lira bani	kina banu sora	0.71	Invalid sentence
moko lira bani	banu kina sora	1.27	Invalid sentence
moko lira bani	sora banu kina	0.75	Invalid sentence
moko lira bani	kina banu tira	1.26	Invalid sentence
moko lira bani	banu tira kina	0.63	Invalid sentence
moko lira bani	tira kina banu	1.46	Invalid sentence
moko lira bani	kina tira banu	0.98	Invalid sentence
moko lira bani	tira banu kina	0.92	Invalid sentence
moko lira bani	banu kina tira	0.95	Invalid sentence
moko lira bani	tira banu lomo	1.60	Invalid sentence
moko lira bani	banu lomo tira	1.84	Invalid sentence
moko lira bani	lomo tira banu	0.83	Invalid sentence
moko lira bani	tira lomo banu	1.42	Invalid sentence
moko lira bani	lomo banu tira	0.87	Invalid sentence
moko lira bani	banu tira lomo	1.87	Valid (count: 1)
koro lumo banu kina	kina sora lumo	16.87	Invalid sentence
moko lira bani	lumo banu sora	0.72	Invalid sentence
moko lira bani	sora lumo banu	1.28	Invalid sentence
moko lira bani	banu lumo sora	0.76	Invalid sentence
moko lira bani	lumo sora banu	1.34	Invalid sentence
moko lira bani	sora banu lumo	0.62	Invalid sentence
moko lira bani	banu sora lumo	1.39	Invalid sentence
moko lira bani	lumo kina sora	0.68	Invalid sentence
moko lira bani	sora lumo kina	1.42	Invalid sentence
moko lira bani	kina lumo sora	0.98	Invalid sentence
moko lira bani	lumo sora kina	1.03	Invalid sentence
moko lira bani	sora kina lumo	0.95	Invalid sentence
moko lira bani	kina sora lumo	1.04	Invalid sentence
moko lira bani	lumo banu tira	0.94	Valid (count: 2)
fanu kina riko	kina riko sora	16.59	Invalid sentence
moko lira bani	riko sora fanu	0.75	Invalid sentence
moko lira bani	sora fanu riko	1.23	Invalid sentence

Oompa Loompa	gemini-2.5-flash	Time (s)	Validity
moko lira bani	fanu riko sora	0.75	Invalid sentence
moko lira bani	riko fanu sora	1.43	Invalid sentence
moko lira bani	sora riko fanu	0.71	Invalid sentence
moko lira bani	fanu sora riko	1.32	Invalid sentence
moko lira bani	kina fanu sora	0.74	Invalid sentence
moko lira bani	fanu sora kina	1.24	Invalid sentence
moko lira bani	sora kina fanu	0.60	Invalid sentence
moko lira bani	kina sora fanu	1.37	Invalid sentence
moko lira bani	sora fanu kina	1.14	Invalid sentence
moko lira bani	fanu kina sora	0.89	Invalid sentence
moko lira bani	riko kina sora	1.09	Invalid sentence
moko lira bani	kina sora riko	1.10	Invalid sentence
moko lira bani	sora riko kina	0.88	Invalid sentence
moko lira bani	riko sora kina	1.09	Invalid sentence
moko lira bani	sora kina riko	0.94	Invalid sentence
moko lira bani	kina riko sora	1.13	Invalid sentence
moko lira bani	fanu riko tira	0.85	Invalid sentence
moko lira bani	riko tira fanu	1.09	Invalid sentence
moko lira bani	tira fanu riko	0.97	Invalid sentence
moko lira bani	fanu tira riko	0.93	Invalid sentence
moko lira bani	tira riko fanu	1.10	Invalid sentence
moko lira bani	riko fanu tira	0.83	Invalid sentence
moko lira bani	fanu riko lumo	1.04	Invalid sentence
moko lira bani	riko lumo fanu	0.87	Invalid sentence
moko lira bani	lumo fanu riko	2.22	Invalid sentence
moko lira bani	fanu lumo riko	0.83	Invalid sentence
moko lira bani	lumo riko fanu	1.32	Invalid sentence
moko lira bani	riko fanu lumo	0.70	Invalid sentence
moko lira bani	fanu lumo tira	1.34	Invalid sentence
moko lira bani	lumo tira fanu	0.89	Valid (count: 3)
koro fanu kina riko	fanu tira lumo	29.51	Invalid sentence
moko lira bani	tira lumo fanu	0.79	Invalid sentence
moko lira bani	lumo fanu tira	1.26	Invalid sentence
moko lira bani	fanu tira lumo	0.78	Invalid sentence
moko lira bani	tira lumo fanu	1.42	Invalid sentence
moko lira bani	lumo fanu tira	0.80	Invalid sentence
moko lira bani	riko tira fanu	0.71	Invalid sentence
moko lira bani	tira fanu riko	0.62	Invalid sentence
moko lira bani	fanu riko tira	0.64	Invalid sentence
moko lira bani	riko tira fanu	1.38	Invalid sentence
moko lira bani	tira fanu riko	0.71	Invalid sentence
moko lira bani	fanu riko tira	2.87	Invalid sentence
moko lira bani	kina riko tira	0.71	Invalid sentence
moko lira bani	riko tira kina	0.58	Invalid sentence

Oompa Loompa	gemini-2.5-flash	Time (s)	Validity
moko lira bani	tira kina riko	0.57	Invalid sentence
moko lira bani	kina tira riko	0.68	Invalid sentence
moko lira bani	tira riko kina	0.92	Invalid sentence
moko lira bani	riko kina tira	0.74	Invalid sentence
moko lira bani	lumo riko kina	1.29	Invalid sentence
moko lira bani	riko kina lumo	0.74	Invalid sentence
moko lira bani	kina lumo riko	0.71	Invalid sentence
moko lira bani	lumo riko kina	0.70	Invalid sentence
moko lira bani	riko kina lumo	0.88	Invalid sentence
moko lira bani	kina lumo riko	0.74	Invalid sentence
moko lira bani	kina sora riko	0.77	Invalid sentence
moko lira bani	sora riko kina	0.84	Invalid sentence
moko lira bani	riko kina sora	0.85	Invalid sentence
moko lira bani	kina sora riko	0.77	Invalid sentence
moko lira bani	sora riko kina	0.79	Invalid sentence
moko lira bani	riko kina sora	1.40	Invalid sentence
moko lira bani	lumo sora kina	0.70	Invalid sentence
moko lira bani	sora kina lumo	0.61	Invalid sentence
moko lira bani	kina lumo sora	0.66	Invalid sentence
moko lira bani	lumo sora kina	0.88	Invalid sentence
moko lira bani	sora kina lumo	1.37	Invalid sentence
moko lira bani	kina lumo sora	0.73	Invalid sentence
moko lira bani	kina tira sora	1.27	Invalid sentence
moko lira bani	tira sora kina	0.75	Invalid sentence
moko lira bani	sora kina tira	0.63	Invalid sentence
moko lira bani	kina tira sora	0.60	Invalid sentence
moko lira bani	tira sora kina	1.42	Invalid sentence
moko lira bani	sora kina tira	0.72	Invalid sentence
moko lira bani	lumo tira sora	0.93	Invalid sentence
moko lira bani	tira sora lumo	0.88	Invalid sentence
moko lira bani	sora lumo tira	0.82	Invalid sentence
moko lira bani	lumo tira sora	0.72	Invalid sentence
moko lira bani	tira sora lumo	0.68	Invalid sentence
moko lira bani	sora lumo tira	0.77	Invalid sentence
moko lira bani	tira lumo kina	0.81	Invalid sentence
moko lira bani	lumo kina tira	0.84	Invalid sentence
moko lira bani	kina tira lumo	0.89	Invalid sentence
moko lira bani	tira lumo kina	0.80	Invalid sentence
moko lira bani	lumo kina tira	0.73	Invalid sentence
moko lira bani	kina tira lumo	0.96	Invalid sentence
moko lira bani	lumo banu tira	0.92	Valid (count: 4)
soro kina batu	soro kina batu	43.01	Valid (count: 5)
koro sanu kina toro	toro sanu kina	0.78	Invalid sentence
moko lira bani	kina toro sanu	0.79	Invalid sentence

Oompa Loompa	gemini-2.5-flash	Time (s)	Validity
moko lira bani	sanu kina toro	0.81	Valid (count: 6)
koro sanu kina tomo	tomo sanu kina	0.66	Invalid sentence
moko lira bani	kina tomo sanu	0.67	Invalid sentence
moko lira bani	sanu kina tomo	0.77	Valid (count: 7)
koro naku tira falu	falu tira naku	0.86	Invalid sentence
moko lira bani	naku falu tira	1.00	Invalid sentence
moko lira bani	tira naku falu	3.00	Invalid sentence
moko lira bani	falu naku tira	1.53	Invalid sentence
moko lira bani	naku tira falu	0.80	Valid (count: 8)
lira tomo fanu	lira tomo fanu	113.67	Valid (count: 9)
koro tika lira fanu	fanu lira tika	0.76	Invalid sentence
moko lira bani	tika fanu lira	0.81	Invalid sentence
moko lira bani	lira tika fanu	0.72	Invalid sentence
moko lira bani	fanu tika lira	0.92	Invalid sentence
moko lira bani	tika lira fanu	0.74	Valid (count: 10)
koro sora nira fanu	fanu nira sora	0.71	Invalid sentence
moko lira bani	sora fanu nira	1.12	Invalid sentence
moko lira bani	nira sora fanu	0.85	Invalid sentence
moko lira bani	fanu sora nira	0.85	Invalid sentence
moko lira bani	sora nira fanu	0.89	Valid (count: 11)
banu tira lomo	lomo banu tira	1.05	Invalid sentence
moko lira bani	tira lomo banu	0.85	Invalid sentence
moko lira bani	banu tira lomo	0.81	Valid (count: 12)
koro lumo banu kina	kina banu lumo	0.83	Invalid sentence
moko lira bani	lumo kina banu	2.01	Invalid sentence
moko lira bani	banu lumo kina	0.85	Invalid sentence
moko lira bani	kina lumo banu	0.79	Invalid sentence
moko lira bani	lumo banu kina	0.84	Valid (count: 13)
koro lumo tira fanu	fanu tira lumo	1.57	Invalid sentence
moko lira bani	lumo fanu tira	0.79	Invalid sentence
moko lira bani	tira lumo fanu	0.85	Invalid sentence
moko lira bani	fanu lumo tira	0.81	Invalid sentence
moko lira bani	lumo tira fanu	0.74	Valid (count: 14)
koro fanu kina riko	riko kina fanu	0.82	Invalid sentence
moko lira bani	fanu riko kina	0.69	Invalid sentence
moko lira bani	kina fanu riko	2.03	Invalid sentence
moko lira bani	riko fanu kina	0.74	Invalid sentence
moko lira bani	fanu kina riko	0.69	Valid (count: 15)
lira banu tomo	lira banu tomo	1.11	Valid (count: 16)
koro sanu lira tomo	tomo lira sanu	0.82	Invalid sentence
moko lira bani	lira tomo sanu	0.75	Invalid sentence
moko lira bani	sanu lira tomo	0.72	Valid (count: 17)
koro sira lira fanu	fanu lira sira	0.80	Invalid sentence
moko lira bani	lira fanu sira	0.79	Invalid sentence

Oompa Loompa	gemini-2.5-flash	Time (s)	Validity
moko lira bani	sira lira fanu	0.92	Valid (count: 18)
koro mako lira sanu	sanu lira mako	4.09	Invalid sentence
moko lira bani	lira sanu mako	0.82	Invalid sentence
moko lira bani	mako lira sanu	0.88	Valid (count: 19)
banu tira lomo	tira lomo banu	1.07	Invalid sentence
moko lira bani	lomo banu tira	0.82	Invalid sentence
moko lira bani	banu tira lomo	1.25	Valid (count: 20)
koro lumo banu kina	kina banu lumo	0.99	Invalid sentence
moko lira bani	lumo kina banu	3.68	Invalid sentence
moko lira bani	banu lumo kina	0.89	Invalid sentence
moko lira bani	kina lumo banu	0.99	Invalid sentence
moko lira bani	lumo banu kina	0.99	Valid (count: 21)
koro lumo tira fanu	fanu tira lumo	0.99	Invalid sentence
moko lira bani	lumo fanu tira	0.86	Invalid sentence
moko lira bani	tira lumo fanu	0.86	Invalid sentence
moko lira bani	fanu lumo tira	0.79	Invalid sentence
moko lira bani	lumo tira fanu	0.72	Valid (count: 25)