DO LARGE LANGUAGE MODELS HAVE LATERAL THINKING IN PUZZLE-SOLVING GAMES?

Anonymous authors

Paper under double-blind review

Abstract

Large Language Models (LLMs) show exceptional skills in a wide range of tasks, with their ability in lateral thinking standing out as a particularly intriguing area. Lateral thinking in LLMs allows them to understand deeper or suggested meanings from the context, which is essential for making sense of complex scenarios, especially in puzzle-solving games. To delve deeper into and improve the lateral thinking capabilities of LLMs in the realm of puzzle-solving, we introduce the "Lateral Thinking Puzzles" and construct the accompanying dataset. Our novel \mathcal{P} uzzle \mathcal{V} erse framework aims to enhance LLMs' lateral thinking in puzzle-solving games. Complementing this, we propose a creativity metric to ensure comprehensive evaluations. Experiments show that the selected LLMs, after being trained with \mathcal{P} uzzle \mathcal{V} erse, have an average improvement of 101.9% compared to their performance before \mathcal{P} uzzle \mathcal{V} erse that trained LLMs perform better in other reasoning tasks.

025 026 027

004

010 011

012

013

014

015

016

017

018

019

021

023

1 Introduction

Lateral thinking, first proposed by De Bono (1970), is a creative problem-solving approach that involves looking at situations from unconventional perspectives to make reasoning. It's quite distinct from logic and often more useful in generating creative and effective solutions. Lateral thinking is contrast with vertical thinking, which is the conventional logical process. While the latter is like digging one hole deeper and deeper, the former requires abandoning the hole and striking off to the sidelines to dig numerous experimental holes.

034 Lateral thinking is important in solving downstream tasks. It encourages us to view problems from various perspectives, leading to more creative solutions. For example, in business management, it helps break traditional thinking patterns, enabling innovative solutions 037 and providing strategic advice that gives companies a competitive edge. In education, cultivating LLMs with lateral thinking abilities allows educators to access tools that foster creative thinking, design 040 engaging learning materials, and encourage students to explore un-041 conventional approaches to problem-solving. In healthcare, lateral 042 thinking can lead to breakthroughs by offering non-traditional diag-043 nostic and treatment suggestions, particularly for rare or complex 044 cases. For instance, Edward Jenner's decision to explore why dairymaids weren't contracting smallpox, instead of why most did, led to 046 the groundbreaking discovery of the smallpox vaccine. Such lateral 047 thinking is also crucial for Large Language Models (LLMs) (Gi-



Figure 1: Different solutions given by a vertical thinker (i.e. LLM) and a lateral thinker (i.e. human), respectively, based on a complex scenario.

adikiaroglou et al., 2024). Xie et al. (2023) emphasize lateral thinking is one of the creative thinking
 process which promote LLMs solve complex problems more effectively. Take the example shown in
 Fig. 1. When facing the complex scenario where a pilot encounters hydraulic system leakage with
 no means of replenishing the fluid, the LLM, such as GPT-4⁻¹, plays the role of vertical thinker that
 provides traditional suggestions, such as contacting air traffic control, etc. However, the solution

⁰⁵³

¹https://chat.openai.com/

given by the human who plays the role as the lateral thinker is using urine which is not a conventional but effective and simple method.

Research on LLMs' lateral thinking in solving downstream tasks is limited. They mainly focus on 057 conventional logical reasoning associated with vertical thinking, which are divided into decomposing tasks and calling external modules. The former includes using Chain of-Thought (CoT) or Auto-CoT to generates reasoning chains (Wei et al., 2022; Zhang et al., 2022), using active learning to stimulate 060 reasoning capabilities (Diao et al., 2023), using a voting strategy to select the most consistent 061 answer output based on different reasoning paths (Wang et al., 2022), etc. The latter includes 062 using frozen LLMs to automatically generate intermediate reasoning steps (Paranjape et al., 2023), 063 decomposing symbolic reasoning, mathematical reasoning, or algorithmic tasks into intermediate 064 steps (Gao et al., 2023), etc. These methods are not enough to make LLMs owe lateral thinking, which necessitates techniques such as challenging assumptions, seeking alternative solutions with 065 analogy, and embracing ambiguity (Xie et al., 2023). 066

067 However, lateral thinking varies across different contexts, making the choice of context for studying 068 lateral thinking an important consideration. For instance, the example mentioned above requires external knowledge or commonsense, such as "Urine is mostly water and can substitute for it in 069 emergencies", while some puzzle-solving games demand creativity and imagination, like the riddle "What kind of dog never bites?" and the answer is "A hot dog". Therefore, in this paper, we 071 choose puzzle-solving games to investigate LLMs' lateral thinking which has two main reasons: i) 072 Puzzle-solving games typically require thinkers to step outside conventional thought patterns and 073 apply creativity and imagination to understand and solve puzzles. ii) These games offer a clear 074 framework and objective that is to find the answer to the puzzle. This makes lateral thinking in 075 puzzle-solving games more quantifiable and researchable compared with other more open-ended or 076 subjective scenarios. 077

To evaluate and enhance LLMs' lateral thinking in puzzle-solving games, we adopt the Lateral Thinking Puzzles (Sloane & MacHale, 1994). Building on the existing lateral thinking puzzles 079 datasets (Jiang et al., 2023; Huang et al., 2023), we construct the largest Lateral Thinking Puzzles dataset (short for "LTP"), which includes riddles, a sequenced set of questions and answers, solutions, 081 and rules. Based on the LTP dataset, we propose \mathcal{P} uzzle \mathcal{V} erse², a baseline framework that improves the lateral thinking in puzzle-solving games of LLMs through assisting them to propose a series of 083 questions to clarify the riddle's solution. In addition, we propose a novel creative metric, including 084 compliance, reasoning, and completeness for evaluating LLMs' lateral thinking capabilities. Accord-085 ing to the experiments, the PuzzleVerse framework can effectively improve LLMs' performance in LTP, resulting in LLMs with advanced lateral thinking in puzzle-solving games. In summary, our 087 study makes three key contributions: i) We construct the largest lateral thinking puzzles dataset. We 088 also propose the creativity metric, adopting it and human metric to evaluate LLMs' lateral thinking in puzzle-solving games. ii) We make an exploration for LLMs' lateral thinking in puzzle-solving 089 games, and then develop a novel PuzzleVerse framework to enhance these capabilities in LLMs. 090 iii) We validate the effectiveness of PuzzleVerse in LLMs' lateral thinking in puzzle-solving games 091 through extensive experiments in LTP dataset and other reasoning tasks. 092

093 094

095

2 Dataset Construction

096 In this section, we construct a novel lateral thinking puzzles dataset (abbreviated as "LTP") for 097 evaluating and enhancing LLMs' lateral thinking capabilities in problem-solving games. Each puzzle 098 in LTP comprises a riddle and its corresponding solution. The solutions for riddles in LTP are generally unconventional. As shown in Fig. 2, the riddle states that recently your mother has been 099 acting strangely, often distracted, and sneaking out at night, and you need to discover the truth. The 100 conventional solution is to suspect that the mother is having an affair or involved in some secret 101 activities. However, the unconventional solution is that the mother is participating in square dancing. 102 She sneaks out at night to practice with the team, and to avoid disturbing others, they all dance silently 103 with headphones on. The final solution that she is involved in square dancing does not reveal any 104 secret or suspicious activities.

²https://anonymous.4open.science/r/haiguitang-EFA7/. We will open-source all data and code after being accepted.

108 Therefore, due to the unconventional nature of the solutions in LTP, LLMs need to employ lateral 109 thinking without relying on traditional reasoning. They are requested to engage in creative and 110 out-of-the-box thinking to arrive at the solution. Since directly providing a solution based on lateral 111 thinking is highly challenging for LLMs, based on the existing lateral thinking puzzles (Sloane & 112 MacHale, 1994), we set the evaluation of LLMs' lateral thinking capabilities in problem-solving games as follows: for a given riddle, an LLM need to employ lateral thinking through asking yes-or-113 no questions to infer the solution. An LLM that can infer the solution with the fewer questions is 114 considered to have stronger lateral thinking capabilities in these problem-solving games. 115

- 116
- Specifically, we initially collect 647 117 Chinese lateral thinking puzzles from 118 various websites like Huiwan³, Baidu Wenku⁴, etc. Utilizing GPT-4, we 119 generate additional puzzles that mir-120 ror the style and structure but have 121 different semantics from the original 122 ones through in-context learning with 123 the prompt in Table 4 (row "RS Gen-124 eration"). After generating new puz-125 zles, to ensure that these data points 126 have not been previously learned by 127 the considered LLMs, we remove the 128 original 647 puzzles and use only the



Figure 2: A representative puzzle, which includes a riddle, its solution, questions, answers, and clues.

129 generated data for LLMs' evaluation and enhancement. To preserve the unique Chinese characteristics of the dataset and account for the significant semantic differences between Chinese and English, we 130 use the collected Chinese data to expand and create a specialized Chinese dataset. This approach 131 ensures the retention of cultural nuances often lost in translation. Each riddle in the generated puzzles 132 includes only the beginning and end of a story, creating a sense of discontinuity. The solutions require 133 unconventional thinking, differing from standard approaches. Each generated puzzle is assessed 134 using GPT-4 to ensure it meets specific criteria, as detailed in Table 3 (row "RS Evaluation"), with 135 each criterion scored as 0 or 1. Puzzles scoring below 3 are discarded, resulting in a final average 136 score of 3.37. 137

138 Subsequently, we employ GPT-4 to create a sequence of questions, answers,

and five supporting clues for each puzzle with the prompt in Table 4 (row 139 "QAC Generation"). The questions strictly adhere to a yes-or-no format and 140 are crafted to incrementally lead to the solution, reflecting the unconventional 141 nature of the puzzles. Items with formatting errors are discarded and regener-142 ated. Clues are designed to hint at the solution but not the exact solution, and 143 answers are confined to "yes," "no," or "irrelevant." Each set of questions, 144 answers, and clues per puzzle is also evaluated with GPT-4 to ensure logical 145 progression without significant leaps, adequately hint at the solution, and 146 correctly answer the questions. The criteria are shown in Table 3 (row "QAC Evaluation"), with each criterion scoring 0 or 1. Similarly, sets scoring below 147 3 are discarded, resulting in a final average score of 3.52. Importantly, for 148 both RS and QAC evaluation, we successively input instructions, such as 149 first asking, "Does the solution require unconventional thinking, differing 150 from standard approaches?" followed by, "Is the overall logic of the puzzle 151 coherent and readable?". This approach migrating the issue where GPT-4, 152 when provided multiple instructions together, may only output partial ratings, 153

Table 1: The statistics of LTP.

Content	Num.
Avg. Tokens (Riddles)	118.4
Max Tokens (Riddles)	200
Min Tokens (Riddles)	50
Avg. Tokens (Solutions)	63.7
Max Tokens (Solutions)	150
Min Tokens (Solutions)	30
Avg. Tokens (Questions)	13.6
Max Tokens (Questions)	25
Min Tokens (Questions)	10
Avg. Tokens (Clues)	4.7
Max Tokens (Clues)	8
Min Tokens (Clues)	2
Avg. Number of Rounds	15.1
Max Number of Rounds	20
Min Number of Rounds	7

such as a single score (e.g., 1) instead of a complete set of scores (e.g., [1,1,1,1,1]).

Finally, we make quality validation to ensure the quality and safety of LTP, even with unavoidable themes like suicide and murder. GPT-4 is used to automatically detect and flag potentially unsafe content, discarding entries with detailed descriptions of violence and horror. This process ensures the dataset maintains its integrity while minimizing potential risks associated with sensitive content to the fullest extent possible. Ultimately, we generate a total of 647,000 distinct puzzles. We then

¹⁶⁰ 161

³https://huiwan.wepie.com/

⁴https://wenku.baidu.com/

Table 2: Comparison	n of other puzzle-related	problem-solving datasets.
The second se	· · · · · · · · · · · · · · · · · · ·	

164	Dataset	Size	Task Type	Language	Difficulty	Evaluation Content	Evaluation Method
165	BRAINTEASER (Jiang et al., 2023)	1,119	Multiple-Choice QA	English	High	Lateral thinking	Model Answering
	LatEval (Huang et al., 2023)	325	Interactive OA	English, Chinese	High	Lateral thinking	Model Asking and Answering
166	Missed Connections (Todd et al., 2024)	250	Puzzle Game	English	Medium to High	Puzzle-solving	Model Answering
	RiddleSense (Lin et al., 2021)	5,700	Multiple-Choice QA	English	High	Commonsense reasoning	Model Answering
167	LTP (Ours)	642,600	Yes-or-No Questions	Chinese	High	Lateral thinking	Model Asking

Table 3: Rating criteria for evaluating puzzles in LTP.

Content	Criteria
RS Evaluation	Does the puzzle contain only the beginning and end of a story, creating a sense of discontinuity? If yes, score 1; otherwise, score 0. Does the solution require unconventional thinking, differing from standard approaches? If yes, score 1; otherwise, score 0. Is the overall logic of the puzzle coherent and readable? If yes, score 1; otherwise, score 0. Does the puzzle contain any overly detailed descriptions of violence or horror? If yes, score -100; otherwise, score 1. (-100 means the puzzle is discarded regardless of other scores if detailed negative descriptions are present.)
QAC Evaluation	 Do the questions strictly adhere to a yes-or-no format? If yes, score 1; otherwise, score 0. Do the questions incrementally lead to the solution with logical coherence and no significant leaps? If yes, score 1; otherwise, score 0. Do the clues hint at but not reveal the solution? If yes, score 1; otherwise, score 0. Are the answers strictly confined to "Yes," "No," or "Irrelevant"? If yes, score 1; otherwise, score 0.

178 select 30% of the entries in LTP for manual rating by three volunteers. The criteria for this manual 179 rating combine the first two sets assessed by GPT-4, as shown in Table 3 (rows "RS Evaluation" 180 and "QAC Evaluation"). Puzzles scoring below 6 are discarded, resulting in a final average score of 6.65 and a final count of 642,600 distinct puzzles. To ensure the reliability and validity of the human ratings, we calculate the Inter-rater Agreement using Krippendorff's Alpha and discard data 182 entries with an agreement lower than 0.7, resulting in a final agreement of 0.83. The statistics of 183 LTP are documented in Table 1 and more samples in LTP are shown in Table 1. We also compare 184 LTP with other puzzle-related problem-solving datasets as shown in Table 2, which suggests that 185 the constructed LTP is currently the largest and most comprehensive dataset especially for lateral thinking puzzles.

3 Methods

189 190 191

192

213 214 215

187 188

162 163

168

181

In this section, we introduce \mathcal{P} uzzle \mathcal{V} erse, a simple framework inspired by ChatGPT⁵ to enhance LLMs' lateral thinking capabilities in puzzle-solving games.

193 Supervised Fine-Tuning. First, we 194 make Supervised Fine-Tuning (SFT) with an LLM. The input consists of riddles, 195 the historical question-answer sequences, 196 and clues with the instruction "Please ask 197 a yes-or-no question based on the riddle [CONTENT], previous question-answer se-199 quences [CONTENT], and clues [CON-200 TENT].", and output the next question. 201 During the training process, we employ 202 scheduled sampling (Bengio et al., 2015) 203 that balances teacher-forcing and free-204 generation. In the initial stages, teacher-205 forcing is used to ensure that the LLM



Figure 3: The overview of *PuzzleVerse* framework.

learns the optimal question generation paths. Questions in the training dataset serve as target 206 ones and are used as input to train the LLM in question generation. As training progresses, free-207 generation is introduced, enabling the LLM to learn to generate questions independently and refine 208 its strategy for progressive questioning. During free-generation, we use the LLM's own generated 209 questions as input and compare these generated questions with the corresponding target question. 210 The proportion of teacher-forcing gradually decreases, and that of free-generation correspondingly 211 increases according to the following equations: 212

$$p = \frac{1}{1 + e^{-\tau(k-k_0)}}, \quad L_s = pL_t + (1-p)L_f, \tag{1}$$

⁵https://chatgpt.com/

218		
219	Content	Prompt
220	RS Generation	Given the following puzzle which contains a riddle [CONTENT] and a solution [CONTENT], generate a new puzzle that mirror the style and structure but have different semantics. The generated puzzle contains the riddle and a solution.
221	QAC Generation	The puzzle is [CONTENT]. Given the puzzle, generate a sequence of yes-or-no questions that incrementally lead to the solution. Then generating an answer of each question. The answers is confined to "yes," "no," or "irrelevant." based on the riddle and the solution. After that, provide five supporting clues that thin at the solution without revealing it directly.
223	Questioning	The riddle is [CONTENT]. [The previous questions and answers are [CONTENT]]. Given the riddle, [the previous questions and answers], please ask a "yes-or-no" question.
224	Answering	Please respond to the question in "Yes" or "No" or "Irrelevant". "Irrelevant" means that the current question is not important to deduce the solution. If the answers to five consecutive questions are either "No" or "Irrelevant", provide a clue from the given clues [CONTENT]. You need to give the sign of [SUCCESS] if the questioner deduces the solution within the round limits. Otherwise, you should give the sign of [FAL].
226	CS Outputting	Given the following riddle [CONTENT], solution [CONTENT], the question [CONTENT], and the answer [CONTENT], please rate the confidence of the answer on a scale of 1 to 5 (1 being the worst and 5 being the best).

216 Table 4: Prompts for data generation, for the interaction between the questioner LLM and answerer LLM in the 217 inference process, and outputting confidence scores.

where p represents the proportion of teacher-forcing, k is the current training step, k_0 is the starting step of the decay, τ is a parameter controlling the decay rate. L_s and L_t represent the respective loss of teacher-forcing and free-generation. s_1 to s_n in Fig. 3(a) represent the states.

Reward Model Construction. Then, we construct a reward model for the generated questions to 231 encourage LLMs to further generate next questions based on the optimal path. Firstly, we adopt 232 GPT-3.5 as the answerer LLM to answer the generated questions with the prompt in Table 4 (row 233 "Answering"). The questions answered as "Yes" receive positive rewards, while the other questions 234 answered as "No" or "Irrelevant" receive negative rewards. And questions answered as "No" have 235 higher rewards than those answered as "Irrelevant". Subsequently, we determine the overlap score 236 between each positive-rewarded question and the solution. The overlap score measures the similarity, 237 evaluated through sentence embedding using SimCSE (Gao et al., 2021), between the question and 238 the solution. Questions with a higher overlap score receive higher rewards. Additionally, we request 239 the answerer LLM to provide a confidence score between 1 and 5 for the generated questions to further refine the rewards. This confidence score reflects the answerer LLM's trust in its own answers, 240 which is inspired by the reliability metric from LLMs' hallucination evaluation metrics proposed 241 by Chen et al. (2023a) with the prompt in Table 4 (row "CS Outputting"). 242

243 We then combine the overlap and confidence scores to compute the reward r_i of a generated question 244 q_i as follows:

$$r_{i} = \begin{cases} \alpha o(q_{i}) + \beta s(a(q_{i})), & \text{if } a(q_{i}) = \text{Yes} \\ -\alpha o(q_{i}) + \beta s(a(q_{i})), & \text{if } a(q_{i}) = \text{No} \\ -\gamma \alpha o(q_{i}) + \beta s(a(q_{i})), & \text{if } a(q_{i}) = \text{Irrelevant} \end{cases}$$
(2)

where $o(q_i)$ and $s(a(q_i))$ represents the overlap score and confidence score by the answerer, respec-249 tively, for question q_i . α and β are hyper-parameters in (0,1), and γ is a hyper-parameter over 1. This 250 process results in a reliably ranked question sequence $\{q_1, q_2, \ldots, q_{k-1}, q_k\}$ from the most irrelevant 251 to the closest to train a reward model. 252

Specifically, we adopt an LLM, substituting the softmax layer with a linear layer, to construct the 253 reward model, which receives a generated question sequence as input and outputs a score indicating 254 the question quality. We form pairwise ranking pairs from the ranking sequence's generated questions 255 and utilize the Pairwise Ranking Loss (Liu et al., 2009) for training as depicted below: 256

$$L_{\theta} = -\frac{1}{\binom{k}{2}} E_{\sim D}[\log(\sigma(r_{\theta}(x, y_w) - r_{\theta}(x, y_l)))], \tag{3}$$

where x represents the original question, y_w and y_l denote the higher-scoring and lower-scoring 260 questions, respectively, in the corresponding ranking pair. r_{θ} represents the scalar output of the reward 261 model, D is the set of ranking pairs, and K denotes the number of generated questions. Through this 262 process, the reward model learns to attribute higher scores (rewards) to superior questions and lower 263 scores (rewards) to inferior questions. 264

Reinforcement Learning. After that, we adopt Reinforcement Learning (RL) based on the reward 265 model to further search the optimal question generation path. The state is defined as the riddles, 266 previous question-answer pairs, and clues, with the action being the next question to ask. We employ 267 the Proximal Policy Optimization (PPO) method (Schulman et al., 2017) for training. 268

227

228

229

230

245 246

247

248

270 **Experiments** 4 271

272 In this section, we select some powerful LLMs to explore their lateral thinking capabilities in puzzle-273 solving games, and further enhance their capabilities with our PuzzleVerse. In this process, an LLM 274 is tasked with formulating questions about a given riddle, then continuing to ask additional questions 275 based on the answers and clues provided by the answerer, who is set to be GPT-3.5.

276 **Experimental Setups.** We conduct our experiments on four Nvidia A100 GPUs, each with 80GB of 277 memory, using PyTorch in Python. For enhanced training efficiency, we utilize DeepSpeed. We set 278 the maximum sequence length for input and output sequences to 1024 and 200 tokens, respectively. 279 The training process is set to 20 epochs. The detailed configuration of the hyperparameters can be 280 found in Table 5. The prompt of questioning during the inference process is shown in Table 4 (row 281 "Questioning").

282 During the inference process, we first adopt GPT-3.5 to generate an answer among "Yes", "No", 283 "Irrelevant" for each posed question. The input is comprised of the riddles, questions, and clues, and 284 the corresponding output is the answers to the questions. Secondly, we adopt GPT-3.5 to determine 285 the optimal moment to provide given clues for the questioner LLM. If a question is asked with a 286 positive answer (i.e., "Yes"), it receives positive score (such as plus 1). Conversely, a negative score 287 (such as minus 1) is assigned for the question. If a series of questions consecutively receives negative 288 scores for more than five rounds, GPT-3.5 is then requested to generate a clue to guide the questioner. 289 Finally, the GPT-3.5 determines the questioning's termination. Questioning terminates either when the questioner LLM successfully infers the solution or when the questioning reaches a predefined 290 round limit (we defined it as 30). We further utilize the threshold of the overlap score, which is set as 291 0.8 tuned through experimentation, to assess the correlation between the sequence of questions and 292 the solution, determining if the solution has been deduced. If the overlap score exceed this threshold 293 within the round limits, it indicates the questioner's successful deduction, prompting GPT-3.5 to 294 declare questioning termination. Alternatively, GPT-3.5 signifies questioning termination when the it 295 reaches the round limits. 296

	Table 5:	Parameter	configuration	and	descriptions.
--	----------	-----------	---------------	-----	---------------

299	Parameter Name	Parameter Value	Parameter Description
200	Teacher Forcing Ratio (p)	0.8	The probcapability of using the actual answer as the next input during training, as opposed to using the model's own prediction.
300	Decay Parameter (τ)	0.9	Rate at which the teacher forcing ratio decreases over time, allowing the model to rely more on its own predictions during training.
000	Decay Start Step (k_0)	1000	The training step at which the decay of the teacher forcing ratio begins.
301	Overlap Score Weight (a)	0.7	Weighting given to the overlap score when determining the relevance of a generated question to its context.
501	Confidence Score Weight (β)	0.3	Weighting given to the confidence score when assessing the quality of a generated question.
202	Penalty for Irrelevant Answer (γ)	-0.2	Deductive value applied when a model-generated answer is deemed irrelevant to the context.
302	PPO Clipping Range (ϵ)	0.2	Hyper-parameter in PPO that prevents the policy update from changing too drastically, ensuring stable training.
202	Policy Loss Weight (μ_2)	0.25	Weight given to the policy loss $L^{clip}(\theta)$ during reinforcement learning training.
303	Value Function Loss Weight (µ3)	0.25	Weight given to the value function loss $L^{VF}(\bar{\theta})$ during reinforcement learning training.
304			

305 Datasets, Baselines and Metrics. LTP is divided into training and validation sets in a 7:3 ratio, with 306 70% of the data used to train LLMs and the remaining 30% used to evaluate the LLMs' performance. 307 Even without training, the same 30% dataset is used for performance evaluation of the LLMs. We 308 also incorporate other reasoning tasks, similar to lateral thinking puzzles, to validate the effectiveness 309 of LLMs trained with PuzzleVerse. These tasks include story datasets (e.g., LOT (Guan et al., 2022)) 310 and reading comprehension datasets (e.g., DuReader (He et al., 2017), MS MARCO (Nguyen et al., 2016)). The evaluation metrics for these datasets remain consistent with those in the original papers: 311 accuracy for story understanding tasks (i.e., ClozeT, SenPos) and BLEU for story generation tasks 312 (i.e., PlotCom, OutGen) and reading comprehension tasks (i.e., DuReader, MS MARCO). 313

314 We choose Baichuan-7B⁶, ChatGLM-6B (Du et al., 2022), BELLE-13B (Yunjie Ji, 2023; Yunjie Ji & 315 Li, 2023), MOSS-16B (Sun et al., 2023), and GPT4 as baseline LLMs to evaluate their lateral thinking 316 capabilities. We also adopt \mathcal{P} uzzle \mathcal{V} erse to enhance the performance of the open-sourced LLMs (the first four LLMs). 317

318 To evaluate the quality of the generated questions, we design a comprehensive set of metrics, including 319 creativity metric, machine metric, and human metric. Creativity metric comprises compliance, reason-320 ing, and completeness scores. Machine metric includes BLEU (Papineni et al., 2002), ROUGE (Lin, 321 2004), the diversity score (Li et al., 2016), and the embedding score (Liu et al., 2016). Human 322 metric is an average score that combines compliance, reasoning, and completeness. Specifically,

297 298 299

³²³

⁶https://github.com/baichuan-inc/Baichuan-7B

Content	Criteria
Creativity Evaluation	Compliance Score: If half or more of the questions in a puzzle are in the yes-or-no format, the score is 1; otherwise, the score is 0. Reasoning Score: If half or more of the follow-up questions in a puzzle are based on previous information, the score is 1; otherwise, the score is 0. Completeness Score: If the correct solution to a puzzle is provided within the limited number of turns, the score is 1; otherwise, the score is 0.
Human Evaluation	If less than half of the questions in a puzzle are in the yes-or-no format, less than half of the follow-up questions are based on previous question-answer p and clues, and the correct solution is not deduced within the limited number of turns, the score is 1. If half of the questions in a puzzle are in the yes-or-no format, half of the follow-up questions are based on previous question-answer pairs and clues, and
	correct solution is not deduced within the limited number of turns, the score is 2. If more than half of the questions in a puzzle are in the yes-or-no format, more than half of the follow-up questions are based on previous question-answer provide the statement of the statement
	and clues, and the correct solution is not deduced within the limited number of turns, the score is 3. If all the questions in a nuzzle are in the vescorno format, all the follow-up questions are based on previous question-answer pairs and clues, and the corr
	solution is not deduced within the limited number of turns, the score is 4.
	If all the questions in a puzzle are in the yes-or-no format, all the follow-up questions are based on previous question-answer pairs and clues, and the cor solution is deduced within the limited number of turns, the score is 5.

Table 7: Frameworks related to lateral thinking capabilities.

Framework	Target Task	Core Technology	Lateral Think- ing Support	Innovation	Performance
Auto-CoT (Zhang et al., 2022)	Logical Reasoning	Automatic Generation of Reasoning Chains	Weak	Traditional reasoning based on logic	Performs well in logical reasoning tasks, but lacks lateral thinking sup- port
PAL (Gao et al., 2023)	Algorithmic Reasoning	Automatic Decomposition of Algo- rithmic Steps	Weak	Focuses on symbolic and algorith- mic reasoning	Performs well in mathematical and algorithmic tasks, but not suitable for lateral thinking
Connections Solver (Todd et al., 2024)	Puzzle Game	Sentence Embeddings and Instruction-Tuned LLMs	Medium	Combines sentence embeddings with LLMs to solve complex puz- zle tasks	Performs well in the "Connections" puzzle task, testing the impact of dif- ferent prompting styles
\mathcal{P} uzzle \mathcal{V} erse (Ours)	Puzzle-Solving and Lateral Thinking	Question Generation and Reasoning Chain Analysis	Strong	Provides novel evaluation metrics	Excels in the LTP dataset

the creativity metric is obtained by GPT-4 to assess how well the LLM adheres to the rules and the effectiveness of its generated questions in achieving the solution with 0-1 scale based on the criteria shown in Table 6 (row "Creativity Evaluation"). Scores in this metric are designed based on the characteristics of the lateral thinking game. For instance, the compliance score evaluates whether the generated questions adhere to the basic rules of yes-or-no answers, a critical element in the game. The reasoning score assesses whether follow-up questions are based on previous question-answer pairs. The strength of reasoning ability directly impacts the progress of the puzzle-solving process, making it a crucial evaluation dimension that reflects whether LLMs possess coherent thinking abilities. The completeness score measures the extent to which the generated questions effectively lead to the solution, directly reflecting the effectiveness of LLMs' lateral thinking. Given that the puzzles are designed to be approached from unconventional angles, questions that systematically lead to the solution are considered crucial for fostering lateral thinking. For human metric, we enlist nine human raters to evaluate questions from 1,000 randomly selected puzzles with a 1-5 scale based on the criteria shown in Table 6 (row "Human Evaluation"). The raters kindly offered their assistance without compensation. Inter-rater agreement, measured using Krippendorff's Alpha, is used to ensure rating confidence. Controversial ratings with low agreement (<0.7) are discarded, and questions from another riddle are selected for evaluation. By combining diverse and comprehensive evaluation, we reduce biases that arise from a single evaluation metric, increasing the reliability and credibility of the scoring.





Table 8: The lateral thinking performance of vanilla LLMs and that of PuzzleVerse-trained LLMs. "PV" means training LLMs with PuzzleVerse.

373																								
374		с	omplianc	e	I	Creativity Reasoning		Co	mpletene	ss		BLEU			ROUGE	Mac	hine I	viversity-2			ES		Hum	ıan
075		w/o PV	w/ AB	w/ PV	w/o PV	w/ AB	w/ PV	w/o PV	w/ AB	w/ PV	w/o PV	w/ AB	w/ PV	w/o PV	w/ AB	w/ PV	w/o PV	w/ AB	w/ PV	w/o PV	w/ AB	w/ PV	w/o PV	w/ PV
375	baichuan	79.5	81.3	84.4	23.4	39.5	57.0	32.3	49.1	68.1	10.9	18.6	31.1	24.3	32.6	43.5	65.8	68.2	72.9	23.5	37.6	55.0	1.9	3.8
070	MOSS	76.0	78.2	84.1	20.5	35.7	56.0	31.4	48.7	67.4	10.3	17.3	30.4	21.0	30.7	42.8	64.3	66.7	72.3	22.3	34.1	54.3	1.7	3.6
376	BELLE	74.7	77.5	83.7	19.6	28.8	48.9	31.1	46.5	51.2	9.7	16.9	30.1	29.7	35.8	49.2	62.1	65.9	72.9	21.0	30.3	53.5	1.4	2.8
	ChatGLM	/2.0	/0.4	83.0	17.8	25.4	40.0	29.5	43.0	51.1	10.0	10.5	30.2	19.9	28.0	40.5	01.3	65	72.8	19.5	28.7	51.1	1.2	2.9
377	Average	75.7	78.4	84.0	20.3	32.4	52.0	31.1	46.8	59.5	10.2	17.3	30.5	23.7	31.4	44.0	63.4	66.5	72.7	21.6	32.7	53.5	1.6	3.3
011	Ť	-	2.7	8.3	-	12.0	31.7	-	15.8	28.4	-	7.1	20.2	-	7.7	20.3	-	3.1	9.3	-	11.1	31.9	1 -	1.7
	†(%)	-	3.5	10.9	-	59.2	155.7	· ·	50.7	91.3	-	69.4	197.8	-	32.5	85.5	-	4.9	14.8	-	51.4	147.9	- 1	111.3

384

386 387 388

389 390 391

Table 9. A comparison of GP1-4 with zero-shot results from other models across 1,000 sample	Tab	le 9: A comparison	of GPT-4 wit	h zero-shot	results from	other models	across 1,00)0 samples
---	-----	--------------------	--------------	-------------	--------------	--------------	-------------	------------

		Creativity		1	Mae	chine		Human
	Compliance	Reasoning	Completeness	BLEU	ROUGE	Diversity-2	ES	/ /
baichuan	77.3	22.6	35.9	11.6	27.5	64.9	24.4	1.9
MOSS	72.4	21.5	33.1	10.1	20.5	64.1	23.5	1.8
BELLE	74.0	18.2	30.6	9.5	28.1	63.8	21.6	1.4
ChatGLM	71.8	17.9	29.3	10.0	21.9	60.3	19.8	1.3
GPT4	91.7	72.5	78.8	56.2	79.3	84.4	70.1	4.3

Table 10: One-shot performance of LLMs in other reasoning tasks after being trained with PuzzleVerse.

	S	tory Und	erstanding			Story Ge	eneration		Reading Comprehension				
	ClozeT		SenPos		PlotCom		OutGen		Dureader		MSMACRO		
	w/o PV	w/ PV	w/o PV	w/ PV	w/o PV	w/ PV	w/o PV	w/ PV	w/o PV	w/ PV	w/o PV	w/ PV	
baichuan	81.7	88.5	70.5	78.4	29.5	34.1	51.2	59.1	49.1	58.3	42.5	47.1	
MOSS	79.3	85.4	67.5	74.6	26.3	30.7	50.4	56.2	47.5	53.9	39.7	45.9	
BELLE	76.9	84.9	68.1	76.8	25.7	30.5	48	55.2	47.2	54.5	38.4	45.7	
ChatGLM	76.1	83.2	64.2	70.8	23.5	27.4	45.2	53.5	46.5	52.3	38.3	44.3	
Average	78.5	85.5	67.6	75.2	26.3	30.7	48.7	56.0	47.6	54.8	39.7	45.8	
↑ ⁻	-	7.0	-	7.6	- 1	4.4	-	7.3	-	7.2	-	6.1	
[↑] (%)	-	8.9	-	11.2	- 1	16.9	-	15.0	-	15.1	-	15.2	

392 Main Results. The lateral thinking performance of vanilla LLMs and that of PuzzleVerse-trained 393 LLMs are shown in Table 8. Results of GPT4 is on 1,000 samples due to resource constraints, and 394 the corresponding zero-shot performance of other baseline LLMs is shown in Table 9. From the initial performance of LLMs (denoted as "w/o PV"), we observe that in compliance, baichuan and 396 MOSS score the highest, while BELLE and ChatGLM score relatively lower. In reasoning, all LLMs 397 score low, with baichuan having the highest score at only 23.4. In completeness, baichuan and MOSS have relatively high scores, whereas other two score lower. Machine metrics show baichuan 399 performing well, while other LLMs also perform similarly overall. In human evaluations, all LLMs have poor performance, with scores not exceeding half. Overall, LLMs' initial lateral thinking 400 401 capabilities are limited, especially in reasoning and completeness. Moreover, we find GPT-4 can better zero-shot solve these puzzles, which serves as a non-trivial reference baseline. After training 402 with PuzzleVerse (denoted as "w/ PV"), all LLMs shows significant improvement, particularly 403 in reasoning and completeness. In compliance, all LLMs improve their scores by approximately 404 10% on average, with the gains being relatively modest due to the high baseline of compliance. 405 The improvement in reasoning is particularly significant, with an average increase of over 150%. 406 Completeness scores and machine metrics also see effective enhancement. In human evaluations, all 407 LLMs show improved scores, with an average increase of over 100%. However, these LLMs still 408 have a long way to go compared with GPT-4. 409

We also compare the performance of *P*uzzle*V*erse-trained LLMs with the agent mentioned in AgentBench (Liu et al., 2023) for the LTP task on our LTP dataset (denoted as "w/ AB"), as shown in Table 8. We adopt both creative metrics and machine metrics for evaluation. We find that *P*uzzle*V*erse achieves better results, with an average improvement of 40.5% over the agent. This improvement is likely because the agent can be considered an external prompt-based method, whereas our approach involves training, which better enhances LLMs' performance. There are also some frameworks related to lateral thinking that are not specifically designed for it. Therefore, we only qualitatively compared their target tasks, core technology, innovation, and performance in Table 7.

In addition, we evaluate PuzzleVerse-trained LLMs on other reasoning tasks, including story understanding, story generation, and reading comprehension, as shown in Table 10. We use a one-shot evaluation method, providing each data point with one example. We find that PuzzleVerse-trained LLMs exhibit significant enhancements compared to vanilla models, highlighting the adaptability of PuzzleVerse across a range of reasoning tasks.

Ablation Study. After that, we adopt an ablation study to evaluate the contributions of each module 423 within the PuzzleVerse framework. Due to the strong correlation between the creativity metric and 424 the human metric, we primarily analyze these two metrics, as highlighted in Fig 4. Detailed results 425 are shown in Tables 11. We observe it is evident that each module within the \mathcal{P} uzzle \mathcal{V} erse framework 426 has a significant impact on lateral thinking. We can see that for all dimensions, the scores decrease 427 when any single module is removed. Notably, removing the teacher-forcing module (denoted as 428 "w/o TF") leads to the largest decline across various dimensions, indicating that the teacher-forcing module plays a crucial role in maintaining overall performance. The next most impactful module 429 is reinforcement learning (denoted as "w/o RL"). Free-generation (denoted as "w/o FG") has the 430 smallest effect across all dimensions, showing minimal decline when removed. For creativity and 431 human evaluations, removing the teacher-forcing module results in substantial decreases in human

	Compliance						Creativity Reasoning				Completeness				Human					
	w/PV	w/o TF	w/o FG	w/o RL	w/o PV	w/PV	w/o TF	w/o FG	w/o RL	w/o PV	w/PV	w/o TF	w/o FG	w/o RL	w/o PV	w/ PV	w/o TF	w/o FG	w/o RL	w/o PV
baichuan	84.4	80.3	83.9	82.6	79.5	57.0	45.8	54.5	52.3	23.4	68.1	50.5	64.9	58.3	32.3	3.8	2.0	3.6	2.9	1.9
MOSS	84.1	78.2	83.5	82.0	76.0	56.0	47.2	55.7	53.5	20.5	67.4	52.7	66.5	62.4	31.4	3.6	2.1	3.3	2.5	1.7
BELLE	83.7	76.2	82.1	80.3	74.7	48.9	32.1	44.1	42.8	19.6	51.2	44.7	49.2	47.2	31.1	2.8	1.9	2.7	2.4	1.4
ChatGLM	83.6	73.5	81.6	77.8	72.6	46.0	31.5	42.6	39.2	17.8	51.1	40.6	49.9	46.8	29.5	2.9	1.5	2.4	2.1	1.2
Average	84.0	77.1	82.8	80.7	75.7	52.0	39.2	49.2	47.0	20.3	59.5	47.1	57.6	53.7	31.1	3.3	1.9	3.0	2.5	1.6
Ļ	-	6.9	1.2	3.3	8.3	-	12.8	2.8	5.0	31.7	-	12.3	1.8	5.8	28.4	-	1.4	0.3	0.8	1.7
↓(%)	- 1	8.2	1.4	3.9	9.8	-	24.7	5.3	9.7	60.9	-	20.7	3.1	9.7	47.7	- 1	42.7	8.4	24.4	52.7

Table 11: Performance of training LLMs with PuzzleVerse variants which are removed a certain module. "w/o TF", "w/o RL", and "w/o FG stand for variants without teacher-forcing, RL, and free-generation, respectively.

scores and reasoning, while compliance sees a smaller decline, likely due to its high baseline. These findings indicate that using the complete PuzzleVerse framework brings the greatest improvement across all metrics, highlighting its positive impact on enhancing LLMs' lateral thinking capabilities in problem-solving games.



Figure 5: A good case of baichuan trained with PuzzleVerse on a lateral thinking puzzle.

Case Study. We analyze a good case as shown in Fig. 5. The LLM first asks about the reason of waking up, ruling out the possibility of thirst or hunger as the reason for waking up. Then it shifts the focus to health and asked, eliminating any association with studying or work. Next, the LLM continues to delve deeper, inquiring, "Do you have a certain disease that requires you to take medication at regular intervals?", further revealing that "sleep forever" has relationship with regular medication. Finally, the LLM asks whether not taking the medication on time threaten your life, confirming that not adhering to the medication schedule could endanger the life and therefore hinting at the cause of demise. Through these precise questions, the LLM successfully deduces that disease in the riddles requires regular medication, and failing to take it on time could be life-threatening. We showcase some bad cases in Table 12. The questions, such as whether the man checks the door lock or call the police, are indeed not directly relevant to the progression of the puzzle. After a series of answers with "Irrelevant", the LLM still asks some irrelevant questions.



Riddle	Solution	Questions and Answers	Clues
"Your takeout has ar- rived." "Okay." Af- ter the door closed, the man closed his eyes in terror. Shortly afterward, he expe- rienced intense fear and anxiety. Let's rea- son this out.	The man lived alone and one night, feeling hungry, he ordered takeout but fell asleep while waiting. In the middle of the night, he was awakened by the sound of some- one knocking on the door with the delivery. As he prepared to get out of bed, he heard someone's voice outside the door and re- alized that there was someone else in his home. Fearfully, he closed his eyes and pretended to be asleep. However, shortly afterward, he heard someone whispering in his ear, saying. 'T know you're not asleep.'	Question: Did the man check the door lock after hearing someone at the door? Answer: Irrelevant. Question: Does the man have surveillance cameras installed in his house? Answer: Not important. Question: Did the man lock the door after closing it? Answer: Irrelevant. Question: Did the man call the police after hearing someone at the door? Answer: Irrelevant. Question: Did the man close his eyes because of psychological fear? Answer: Irrelevant.	Takeout delivery Late at night Sounds at the door Closing the door Someone whispering in your ear

Related Work

Puzzle Solving. For example, Jiang et al. (2023) introduced a multiple-choice QA task designed to test and benchmark the lateral thinking abilities of LLMs. Huang et al. (2023) proposed LatE

486 Zhao & Anderson (2023) focused on the ability of LLMs to solve and create puzzles in NPR Sunday 487 Puzzles. King (2023) pointed out the challenges LLMs face in generating anagrams. Zhang et al. 488 (2024) introduced a novel solver-layer adaptation (SoLA) method that enhances the puzzle-solving 489 capabilities of LLMs. Wu et al. (2023) delved into the use of GPT-4 for tackling more complex 490 mathematical problems. Xie et al. (2023) proposed OlaGPT to approximate various cognitive processes, including reasoning and decision-making. Sarathy et al. (2024) introduced ESCAPE using 491 puzzle video games to study cognitive processes in creative problem-solving. Wang et al. (2024) 492 player behavior in a puzzle game to identify effective problem-solving strategies. Differently, our 493 research explore the potential of LLMs in lateral thinking within puzzle-solving games. 494

495 Although some work focus on lateral thinking puzzles and their application in evaluating LLMs, they 496 only provides evaluations without offering solutions. For example, Jiang et al. (2023) introduced a multiple-choice QA task designed to test and benchmark the lateral thinking abilities of LLMs. Huang 497 et al. (2023) proposed LatEval, an interactive benchmark that challenged LLMs on lateral thinking by 498 assessing the quality of questions posed and the integration of information during problem-solving. 499 Todd et al. (2024) explored the use of the "Connections" puzzle game as a benchmark for evaluating 500 LLMs' abstract reasoning and semantic understanding. León Corrales et al. (2010) investigated how 501 lateral thinking puzzles could enhance critical thinking and motivation in students' opinion paragraph 502 writing, leading to improved writing skills. Lin et al. (2021) introduced a multiple-choice QA task focused on riddle-style questions that required commonsense reasoning and linguistic creativity, with 504 a dataset of 5.7k examples. In contrast to these methods, we use LLMs for supervised fine-tuning 505 and reinforcement learning, dynamically generating and optimizing question-posing paths, which 506 significantly improved model performance on LTP tasks. Moreover, none of these benchmarks has as 507 many samples as our work.

Reasoning. For example, Hao et al. (2023) utilized LLMs as world state predictors and strategic reasoners. Lu et al. (2023) introduced Chameleon in enhancing LLMs' compositional reasoning capability. Tarau (2023) automated deep reasoning in LLM dialog threads. Kiciman et al. (2023) delved into causal reasoning capabilities of LLMs. Yoneda et al. (2023) introduced Statler to enhance LLMs' long-horizon reasoning capability in robotic tasks. Paranjape et al. (2023) presented ART to generate intermediate reasoning steps. Chen et al. (2023c) introduced ChatCoT by chain-of-thought reasoning. However, these work mainly focus on vertical thinking instead of lateral thinking.

Question Generation. For example, Chen et al. (2019) designed a reinforcement learning model for natural question generation. Tavares et al. (2023) delved into LLM strategies in generating questions on dialogue state tracking. Kai et al. (2021) proposed a double-hints method for visual question generation. Uehara et al. (2022) stressed the significance of sub-questions in enhancing primary visual queries. Arora et al. (2022) explored effective prompting strategies for LLMs. Abdelghani et al. (2022) harness GPT-3's capabilities in children's curiosity-driven questioning. However, these studies focus on reshaping question generation instead of searching valuable questioning points.

522 Story Understanding. For example, Yuan et al. (2022) introduced a platform fostering human-LLM 523 story-writing collaborations. Swanson et al. (2021) unveiled STORY CENTAUR, optimizing LLMs 524 for creative endeavors. Dong et al. (2022) spotlighted CoRRPUS to boost story consistency in LLM 525 outputs. Bhandari & Brennan (2023) assessed the trustworthiness of LLM-generated children's stories. Chen et al. (2023b) advocated for LLMs to generate complex narratives. Lee et al. (2022) 526 explored LLM-enabled interactive story rewriting. Méndez & Gervás (2023) utilized ChatGPT in 527 narrative "sifting." Together, these contributions highlight the potential of LLMs in story generation 528 and comprehension. 529

530 531

532

6 Conclusions and Future Work

In exploring the potential of LLMs, we've pinpointed their impressive aptitude for lateral thinking,
which is instrumental for grasping intricate and nuanced contexts. By introducing the Lateral
Thinking Puzzles and its complementary dataset, we illuminate the depth of this capability within
LLMs. Our proposed PuzzleVerse framework is designed to further enhance LLMs' lateral thinking
capabilities, and our proposed creativity metric offers a comprehensive evaluation. Experiments show
the effectiveness of PuzzleVerse in not only LTP but also other reasoning tasks. Future research can
delve into more intricate thinking scenarios and introduce the integration of multi-modal data, further
enhancing LLMs' lateral thinking in puzzle-solving games.

540 Ethic Statement

541 542 543

544

546

547

548

549

550

We analyze potential negative impacts and make ethic statement. Firstly, although lateral thinking encourages creativity and non-traditional solutions, these solutions may not align with societal norms or ethical standards in practical applications. Secondly, enhancing lateral thinking capabilities might exacerbate existing biases in LLMs. The previous training data for LLMs may already contain societal biases, and in lateral thinking tasks, these biases could be amplified or perpetuated through the generation of non-traditional solutions. To address these issues, we conduct a more comprehensive analysis of the societal impacts of these capabilities and explore how to incorporate stricter bias detection and correction mechanisms in model development and evaluation. Additionally, ethical reviews are integrated into the evaluation framework of model applications to ensure that the enhancement of lateral thinking capabilities does not lead to adverse societal consequences.

551 552 553

554

Reproducibility Statement

Part of source code is available in https://anonymous.4open.science/r/haiguitang-EFA7/. We will
 open-source all data and code after being accepted. We make reproducibility statement on data
 construction as follows:

558 **Dataset Composition.** We constructed a novel lateral thinking puzzles dataset (LTP) to evaluate 559 and enhance LLMs' lateral thinking capabilities in problem-solving games. Each puzzle includes 560 a riddle with an unconventional solution, requiring creative, out-of-the-box thinking. We initially collected 647 Chinese lateral thinking puzzles from websites like Huiwan and used GPT-4 to generate 561 additional puzzles with different semantics. These were carefully curated and expanded to maintain 562 cultural nuances, resulting in a final dataset of 642,600 puzzles. Each puzzle includes questions, 563 answers, and clues to guide LLMs towards the solution, evaluated for logical progression and safety. The comprehensive LTP dataset offers a robust framework for assessing and improving LLMs' lateral 565 thinking abilities. 566

567 **Collection Process.** We constructed the Lateral Thinking Puzzles (LTP) dataset to enhance and evaluate LLMs' lateral thinking capabilities. Initially, we gathered 647 Chinese puzzles from 568 websites like Huiwan. Using GPT-4, we generated additional puzzles with different semantics to 569 ensure originality. Each puzzle includes a riddle and an unconventional solution, requiring creative 570 thinking beyond traditional reasoning. To preserve cultural nuances, we focused on expanding the 571 dataset in Chinese. We used GPT-4 to create sequences of yes-or-no questions, answers, and clues 572 for each puzzle, designed to guide LLMs toward the solution. Both the puzzles and the question 573 sequences were rigorously evaluated to ensure logical consistency and quality. To ensure safety, 574 we filtered out puzzles with potentially harmful content. This meticulous process resulted in a 575 high-quality dataset of 642,600 puzzles, providing a robust tool for assessing and improving the 576 lateral thinking capabilities of LLMs in problem-solving games.

577 **Preprocessing/cleaning/labeling.** To ensure the quality and safety of the LTP dataset, we imple-578 mented a thorough preprocessing, cleaning, and labeling process. Initially, we used GPT-4 to generate 579 additional puzzles, ensuring they mirrored the style of the collected Chinese puzzles but with different 580 semantics. Each generated puzzle underwent rigorous evaluation to meet specific criteria, such as 581 logical consistency and cultural relevance. Puzzles scoring below a threshold were discarded. Next, 582 we created sequences of yes-or-no questions, answers, and clues for each puzzle, designed to guide 583 the LLMs incrementally towards the solution. These sequences were evaluated for logical progression 584 and accuracy, with inadequate sets being discarded. To maintain dataset integrity and minimize risks, we used GPT-4 to automatically detect and flag potentially unsafe content, such as detailed 585 descriptions of violence or horror. Entries containing such content were removed. Manual rating 586 by volunteers further ensured the dataset's quality, with puzzles scoring below a set threshold being 587 excluded. The final dataset, comprising 642,600 puzzles, was thoroughly vetted for reliability and 588 cultural nuance, ensuring it serves as a robust tool for enhancing LLMs' lateral thinking capabilities. 589

References

591 592

590

Rania Abdelghani, Yen-Hsiang Wang, Xingdi Yuan, Tong Wang, Hélène Sauzéon, and Pierre-Yves Oudeyer. Gpt-3-driven pedagogical agents for training children's curious question-asking skills.

594 505	arXiv preprint arXiv:2211.14228, 2022.
595 596 597 598	Simran Arora, Avanika Narayan, Mayee F Chen, Laurel J Orr, Neel Guha, Kush Bhatia, Ines Chami, Frederic Sala, and Christopher Ré. Ask me anything: A simple strategy for prompting language models. <i>arXiv preprint arXiv:2210.02441</i> , 2022.
599 600 601	Samy Bengio, Oriol Vinyals, Navdeep Jaitly, and Noam Shazeer. Scheduled sampling for sequence prediction with recurrent neural networks. <i>Advances in neural information processing systems</i> , 28, 2015.
602 603 604	Prabin Bhandari and Hannah Marie Brennan. Trustworthiness of children stories generated by large language models. <i>arXiv preprint arXiv:2308.00073</i> , 2023.
605 606	Yu Chen, Lingfei Wu, and Mohammed J Zaki. Reinforcement learning based graph-to-sequence model for natural question generation. <i>arXiv preprint arXiv:1908.04942</i> , 2019.
608 609 610 611	Yuyan Chen, Qiang Fu, Yichen Yuan, Zhihao Wen, Ge Fan, Dayiheng Liu, Dongmei Zhang, Zhixu Li, and Yanghua Xiao. Hallucination detection: Robustly discerning reliable answers in large language models. In <i>Proceedings of the 32nd ACM International Conference on Information and Knowledge Management</i> , pp. 245–255, 2023a.
612 613	Zexin Chen, Eric Zhou, Kenneth Eaton, Xiangyu Peng, and Mark Riedl. Ambient adventures: Teaching chatgpt on developing complex stories. <i>arXiv preprint arXiv:2308.01734</i> , 2023b.
614 615 616 617	Zhipeng Chen, Kun Zhou, Beichen Zhang, Zheng Gong, Wayne Xin Zhao, and Ji-Rong Wen. Chatcot: Tool-augmented chain-of-thought reasoning on\\chat-based large language models. <i>arXiv preprint</i> <i>arXiv:2305.14323</i> , 2023c.
618	Edward De Bono. Lateral thinking. New York, pp. 70, 1970.
619 620 621	Shizhe Diao, Pengcheng Wang, Yong Lin, and Tong Zhang. Active prompting with chain-of-thought for large language models. <i>arXiv preprint arXiv:2302.12246</i> , 2023.
622 623 624	Yijiang River Dong, Lara J Martin, and Chris Callison-Burch. Corrpus: Detecting story incon- sistencies via codex-bootstrapped neurosymbolic reasoning. <i>arXiv preprint arXiv:2212.10754</i> , 2022.
625 626 627 628	Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. Glm: General language model pretraining with autoregressive blank infilling. In <i>Proceedings of the 60th</i> <i>Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pp. 320–335, 2022.
630 631 632	Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. Pal: Program-aided language models. In <i>International Conference on Machine Learning</i> , pp. 10764–10799. PMLR, 2023.
633 634	Tianyu Gao, Xingcheng Yao, and Danqi Chen. Simcse: Simple contrastive learning of sentence embeddings. <i>arXiv preprint arXiv:2104.08821</i> , 2021.
635 636 637 638	Panagiotis Giadikiaroglou, Maria Lymperaiou, Giorgos Filandrianos, and Giorgos Stamou. Puzzle solving using reasoning of large language models: A survey. <i>arXiv preprint arXiv:2402.11291</i> , 2024.
639 640 641	Jian Guan, Zhuoer Feng, Yamei Chen, Ruilin He, Xiaoxi Mao, Changjie Fan, and Minlie Huang. Lot: A story-centric benchmark for evaluating chinese long text understanding and generation. <i>Transactions of the Association for Computational Linguistics</i> , 10:434–451, 2022.
642 643 644 645	Shibo Hao, Yi Gu, Haodi Ma, Joshua Jiahua Hong, Zhen Wang, Daisy Zhe Wang, and Zhiting Hu. Reasoning with language model is planning with world model. <i>arXiv preprint arXiv:2305.14992</i> , 2023.
646 647	Wei He, Kai Liu, Jing Liu, Yajuan Lyu, Shiqi Zhao, Xinyan Xiao, Yuan Liu, Yizhong Wang, Hua Wu, Qiaoqiao She, et al. Dureader: a chinese machine reading comprehension dataset from real-world applications. <i>arXiv preprint arXiv:1711.05073</i> , 2017.

648 Shulin Huang, Shirong Ma, Yinghui Li, Mengzuo Huang, Wuhe Zou, Weidong Zhang, and Hai-Tao 649 Zheng. Lateval: An interactive llms evaluation benchmark with incomplete information from 650 lateral thinking puzzles. arXiv preprint arXiv:2308.10855, 2023. 651 Yifan Jiang, Filip Ilievski, and Kaixin Ma. Brainteaser: Lateral thinking puzzles for large language 652 model. arXiv preprint arXiv:2310.05057, 2023. 653 654 Shen Kai, Lingfei Wu, Siliang Tang, Yueting Zhuang, Zhuoye Ding, Yun Xiao, Bo Long, et al. 655 Learning to generate visual questions with noisy supervision. Advances in Neural Information 656 Processing Systems, 34:11604–11617, 2021. 657 Emre Kıcıman, Robert Ness, Amit Sharma, and Chenhao Tan. Causal reasoning and large language 658 models: Opening a new frontier for causality. arXiv preprint arXiv:2305.00050, 2023. 659 660 Michael King. Large language models are extremely bad at creating anagrams. 2023. 661 Yoonjoo Lee, Tae Soo Kim, Minsuk Chang, and Juho Kim. Interactive children's story rewriting 662 through parent-children interaction. In Proceedings of the First Workshop on Intelligent and 663 Interactive Writing Assistants (In2Writing 2022), pp. 62–71, 2022. 664 665 Helga Valeska León Corrales et al. The use of lateral thinking puzzles to improve opinion paragraph writing.: thinking puzzles to unpuzzle thinking. Master's thesis, Universidad de La Sabana, 2010. 666 667 Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. A diversity-promoting 668 objective function for neural conversation models. In Proceedings of the 2016 Conference of the 669 North American Chapter of the Association for Computational Linguistics: Human Language 670 Technologies, pp. 110–119, San Diego, California, 2016. Association for Computational Linguistics. 671 doi: 10.18653/v1/N16-1014. URL https://www.aclweb.org/anthology/N16-1014. 672 Bill Yuchen Lin, Ziyi Wu, Yichi Yang, Dong-Ho Lee, and Xiang Ren. Riddlesense: Reasoning 673 about riddle questions featuring linguistic creativity and commonsense knowledge. arXiv preprint 674 arXiv:2101.00376, 2021. 675 676 Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In Text Summarization Branches Out, pp. 74-81, Barcelona, Spain, 2004. Association for Computational Linguistics. 677 678 URL https://www.aclweb.org/anthology/W04-1013. 679 Chia-Wei Liu, Ryan Lowe, Iulian Serban, Mike Noseworthy, Laurent Charlin, and Joelle Pineau. How 680 NOT to evaluate your dialogue system: An empirical study of unsupervised evaluation metrics for 681 dialogue response generation. In Proceedings of the 2016 Conference on Empirical Methods in 682 Natural Language Processing, pp. 2122–2132, Austin, Texas, 2016. Association for Computational 683 Linguistics. doi: 10.18653/v1/D16-1230. URL https://www.aclweb.org/anthology/D16-1230. 684 Tie-Yan Liu et al. Learning to rank for information retrieval. Foundations and Trends® in Information 685 Retrieval, 3(3):225-331, 2009. 686 687 Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, 688 Kaiwen Men, Kejuan Yang, et al. Agentbench: Evaluating llms as agents. arXiv preprint 689 arXiv:2308.03688, 2023. 690 Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, 691 and Jianfeng Gao. Chameleon: Plug-and-play compositional reasoning with large language models. 692 arXiv preprint arXiv:2304.09842, 2023. 693 694 Gonzalo Méndez and Pablo Gervás. Using chatgpt for story sifting in narrative generation. 2023. 695 Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, and 696 Li Deng. Ms marco: A human-generated machine reading comprehension dataset. 2016. 697 Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In Proceedings of the 40th Annual Meeting of the Association 699 for Computational Linguistics, pp. 311–318, Philadelphia, Pennsylvania, USA, 2002. Association 700 for Computational Linguistics. doi: 10.3115/1073083.1073135. URL https://www.aclweb.org/ 701 anthology/P02-1040.

702 703 704	Bhargavi Paranjape, Scott Lundberg, Sameer Singh, Hannaneh Hajishirzi, Luke Zettlemoyer, and Marco Tulio Ribeiro. Art: Automatic multi-step reasoning and tool-use for large language models. <i>arXiv preprint arXiv:2303.09014</i> , 2023.
705 706 707 708	Vasanth Sarathy, Nicholas Rabb, Daniel M Kasenberg, and Matthias Scheutz. Using puzzle video games to study cognitive processes in human insight and creative problem-solving. In <i>Proceedings of the Annual Meeting of the Cognitive Science Society</i> , volume 46, 2024.
709 710 711	John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. <i>arXiv preprint arXiv:1707.06347</i> , 2017.
712 713 714	Paul Sloane and Des MacHale. <i>Great lateral thinking puzzles</i> . Sterling Publishing Company, Inc., 1994.
715 716 717 718	Tianxiang Sun, Xiaotian Zhang, Zhengfu He, Peng Li, Qinyuan Cheng, Hang Yan, Xiangyang Liu, Yunfan Shao, Qiong Tang, Xingjian Zhao, Ke Chen, Yining Zheng, Zhejian Zhou, Ruixiao Li, Jun Zhan, Yunhua Zhou, Linyang Li, Xiaogui Yang, Lingling Wu, Zhangyue Yin, Xuanjing Huang, and Xipeng Qiu. Moss: Training conversational language models from synthetic data. 2023.
719 720 721 722	Ben Swanson, Kory Mathewson, Ben Pietrzak, Sherol Chen, and Monica Dinalescu. Story centaur: Large language model few shot learning as a creative writing tool. In <i>Proceedings of the 16th</i> <i>Conference of the European Chapter of the Association for Computational Linguistics: System</i> <i>Demonstrations</i> , pp. 244–256, 2021.
723 724 725	Paul Tarau. Full automation of goal-driven llm dialog threads with and-or recursors and refiner oracles. <i>arXiv preprint arXiv:2306.14077</i> , 2023.
726 727 728 729	Diogo Tavares, David Semedo, Alexander Rudnicky, and Joao Magalhaes. Learning to ask questions for zero-shot dialogue state tracking. In <i>Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , pp. 2118–2122, 2023.
730 731	Graham Todd, Tim Merino, Sam Earle, and Julian Togelius. Missed connections: Lateral thinking puzzles for large language models. <i>arXiv preprint arXiv:2404.11730</i> , 2024.
732 733 734 735	Kohei Uehara, Nan Duan, and Tatsuya Harada. Learning to ask informative sub-questions for visual question answering. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 4681–4690, 2022.
736 737 738	Karen D Wang, Haoyu Liu, David DeLiema, Nick Haber, and Shima Salehi. Discovering play- ers' problem-solving behavioral characteristics in a puzzle game through sequence mining. In <i>Proceedings of the 14th Learning Analytics and Knowledge Conference</i> , pp. 498–506, 2024.
739 740 741 742	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdh- ery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. <i>arXiv preprint arXiv:2203.11171</i> , 2022.
743 744 745	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in Neural Information Processing Systems</i> , 35:24824–24837, 2022.
746 747 748 749	Yiran Wu, Feiran Jia, Shaokun Zhang, Qingyun Wu, Hangyu Li, Erkang Zhu, Yue Wang, Yin Tat Lee, Richard Peng, and Chi Wang. An empirical study on challenging math problem solving with gpt-4. <i>arXiv preprint arXiv:2306.01337</i> , 2023.
750 751 752 753	Yuanzhen Xie, Tao Xie, Mingxiong Lin, WenTao Wei, Chenglin Li, Beibei Kong, Lei Chen, Chengx- iang Zhuo, Bo Hu, and Zang Li. Olagpt: Empowering llms with human-like problem-solving abilities. <i>arXiv preprint arXiv:2305.16334</i> , 2023.
754 755	Takuma Yoneda, Jiading Fang, Peng Li, Huanyu Zhang, Tianchong Jiang, Shengjie Lin, Ben Picker, David Yunis, Hongyuan Mei, and Matthew R Walter. Statler: State-maintaining language models for embodied reasoning. arXiv preprint arXiv:2306.17840, 2023.

- Ann Yuan, Andy Coenen, Emily Reif, and Daphne Ippolito. Wordcraft: story writing with large language models. In 27th International Conference on Intelligent User Interfaces, pp. 841–852, 2022.
- Yan Gong Yiping Peng Qiang Niu Baochang Ma Yunjie Ji, Yong Deng and Xiangang Li. Belle: Be
 everyone's large language model engine. https://github.com/LianjiaTech/BELLE, 2023.
- Yan Gong Yiping Peng Qiang Niu Lei Zhang Baochang Ma Xiangang Li Yunjie Ji, Yong Deng.
 Exploring the impact of instruction data scaling on large language models: An empirical study on real-world use cases. *arXiv preprint arXiv:2303.14742*, 2023.
 - Yu Zhang, Hui-Ling Zhen, Zehua Pei, Yingzhao Lian, Lihao Yin, Mingxuan Yuan, and Bei Yu. Sola: Solver-layer adaption of llm for better logic reasoning. *arXiv preprint arXiv:2402.11903*, 2024.
 - Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. Automatic chain of thought prompting in large language models. *arXiv preprint arXiv:2210.03493*, 2022.
 - Jingmiao Zhao and Carolyn Jane Anderson. Solving and generating npr sunday puzzles with large language models. *arXiv preprint arXiv:2306.12255*, 2023.
- 771 772 773 774

765

766

767 768

769

770

Table 13: Sample puzzles in LTP including riddles, solutions, questions, answers and clues.

775	Riddle	Solution	Questions and Answers	Clues
776	One night, Xiao Ming made a phone call but	Xiao Ming was resting in a hotel when he couldn't fall asleep due to the loud snoring coming from the	Question: Is the phone the one Xiao Ming uses to make calls to others? Answer: Yes. Question: Is the phone the one Xiao Ming uses to call the hotel reception? Answer: Not important.	Resting in a hotel Snoring disrupts
777	it was answered with-	to wake up the person next door, and quickly fell asleen while they were no longer snoring	Question: Is the phone the one Xiao Ming uses to cart the negmoting room: Answer, res. Question: Is the phone the one Xiao Ming internionally hangy up on the other person? Answer: Yes. Question: Is the phone the one Xiao Ming uses to create noise? Answer: Yes. Question: Is the phone the one Xiao Ming uses to complain about the noise to the neighboring room Answer: Not important.	bor Unable to fall asleen
778	other person to speak. Why?	useep while dey were no longer shoring.		chable to har asteep
779			Question: Is the phone the one Xiao Ming uses to request the neighboring room to quiet down? Answer: Not important.	
780			Question: Is the phone the one Xiao Ming uses to communicate with the neighboring room? Answer: Not important.	
781			Question: Is the phone the one Xiao Ming uses to disturb the sleep of the neighboring room? Answer: Yes. Question: Is the phone the one Xiao Ming uses to remind the neighboring room to stop snoring? Answer: Yes.	
782	The woman went to the library to borrow	The woman is the author of this book. She inserted 100 yuan into the book and donated it to the library.	Question: Is the book written by the woman herself? Answer: Yes. Question: Is the book borrowed by the woman from the library? Answer: Yes.	The woman borrowed a book
783	a book. When she opened it, she cried.	After several years, she came back to borrow the same book, only to find the money still inside. This	Question: Is the book the one the woman borrowed from the library? Answer: Yes. Question: Is there a currency note in the book? Answer: Yes.	she cried There was money in-
784		indicates that nobody has actually read her book.	Question: Why did the woman cry? Answer: Not important. Question: How much money did the woman put between the pages of the book? Answer: Not important.	side the book She donated it to the
785			Question: Due the woman donate the book with the money inside to the inerary? Answer? res. Question: How long did it take for the woman to come back to borrow the book? Answer: Not important. Question: Is the money still inside the book? Answer: Yes.	The money is still in- side the book.
786	In a tall building at	The woman saw an ongoing murder incident in	Question: Does the situation imply that hobody looked at the woman's book? Answer: res. Question: Is the woman hanging clothes out at night? Answer: Yes.	At night
787	night, a woman was hanging clothes on the balcony. Sud-	the building across from hers, and the murderer also noticed her witnessing the event. The reason the woman was instantly horrified was that the	Question: Is the woman in a high-rise building where she lives? Answer: Yes. Question: Is the woman hanging clothes on the balcon? Answer: Yes. Ouestion: Did the woman accidentally look towards the building across the street? Answer: Yes.	The woman looked towards the building
788	denly, she uninten- tionally glanced at the	murderer was counting the number of floors in her building.	Joors in her Question: Did the woman see something happening in the building across the street? Answer: Yes. Question: Did the woman witness an ongoing murder incident? Answer: Yes.	
789	building across from hers and was instantly		Question: Did the killer notice that the woman witnessed his actions? Answer: Yes. Question: Did the woman feel terrified because she realized she had been discovered? Answer: Yes.	ing sensation The murderer was
790	horrified.		Question: Is the killer counting the number of floors where the woman is located? Answer: Yes. Question: Does the number of floors where the woman is located have significance to the killer? Answer: Yes.	counting the number of floors.
791	A wealthy man made	In the wealthy man's house, a burglar entered. While the wealthy man was making a phone call	Question: Was the wife at home when the millionaire called her? Answer: Not important. Question: What was the reason for the millionaire to call his wife? Answer: Not important	The wealthy man called his wife
792 793	beloved wife, and as a result, she died.	his wife was hiding in a certain place. Due to the phone not being on silent mode, the ringtone sounded and exposed the wife's location, leading	Question: Is the phone the one the millionaire used to call his wife? Answer: Yes. Question: Did a thie? enter the millionaire's house? Answer: Yes. Question: Did a thie? enter the millionaire's house? Answer: Yes. Question: Did the wife's location get exposed after the phone rang? Answer: Yes. Question: Did the wife's location get exposed after the phone rang? Answer: Yes. Question: Was the wife killed begues of the ringing of the phone? Answer: Yes.	His wife died A burglar entered the house The phone's ringtone sounded The wife's location
794		to her being killed by the burglar.		
795			Question: Did the thief kill the wife because he knew her location? Answer: Yes. Question: Did the thief kill the wife after discovering her hiding place? Answer: Yes.	was exposed.
796	The painter received a phone call, and as he	The painter is a single father, and because his son constantly asked about his mother, he told his son	Question: Is the painter single? Answer: Not important. Ouestion: Is the phone call the painter received an important event? Answer: Yes	The painter received a phone call
797	looked at a mermaid painting on the table,	that the mother is the mermaid in the painting. The young son took it seriously and always said he	Question: Did the painter create the mermaid painting he saw? Answer: Yes. Question: Does the painter's son believe that his mother is the mermaid in the painting? Answer: Yes. Question: Was the painter's son sent to a mental hospital because he was searching for his mother? Answer Yes.	There was a mermaid painting His son was sent to a mental hospital His son died by sui- cide drowning
798	he suddenly started crying.	y started wanted to go into the water to find his mother. Due to this situation, he was eventually sent to a mental hospital. The painter received a call from the mental hospital. The painter received a call from the mental hospital, informing him of his son's suicide by the painter's son die from a suicide by the painter		
799			Question: Is the phone call the painter received about his son? Answer: Yes. Question: Did the painter's son die from a suicide by drowning? Answer: Yes. Ouestion: Does the painter search tendinistic his con's month lieurae orthoging? Answer: Yes.	
800		painting on the table, he deeply regreted not realize ing his son's mental issues earlier or explaining cleases the painter regret not explaining clearly about his son's mother? Answer: Yes.		gretted his past ac- tions.
801		situation clearly, which ultimately led to his son's tragic suicide.		
802	That painting de- picted a man with sharp features, vividly lifelike. The next day, when I saw the painting again. I felt a tingling	I entered a rundown small hotel late at night. When I entered the room, even the light was broken, and	Question: Is the painting in a run-down small hotel? Answer: Yes. Question: Is the man in the painting very handsome? Answer: Yes.	Rundown small hotel Man with sharp fea- tures Feeling uncomfort- able Window The owner mistook it for a painting.
803		the room was dimly lit. There was a painting on the opposite side of the bed, depicting a man with	Question: Is the man in the painting depicted with clear features and lifelike appearance? Answer: Yes, Question: Does the man in the painting make the owner uncomfortable? Answer: Yes, Question: Was the painting latter discovered to be a window by the owner? Answer: Yes, Question: Was the painting latter discovered to be a window by the owner? Answer: Yes, Question: Was the location of the window mistaken for a painting by the owner? Answer: Yes, Question: Did the owner feel that the lighting was dim when looking at the window at night? Answer: Yes Question: Did the owner mistake the man standing outside the window for a painting? Answer: Yes, Question: Did the owner mistake that man on the window was continuously watching him? Answer: Yes	
804		sharp features, vividly lifelike, just like the Mona Lisa. I always felt that the person in the painting was constantly watching me. It wasn't until the		
805	sensation on my scalp, and I couldn't	next morning, when it was bright outside, that I realized the supposed painting was actually a win-		
806	utter a single word of praise.	dow. Last night, there was a man standing outside the window watching me, but due to the dim light, I mistook him and the window frame for a point	Question: Did the owner only discover that it was actually a window the next morning? Answer: Yes.	
807		ing.		
000				

808