



# Evolutionary System 2 Reasoning: An Empirical Proof

Zeyuan Ma<sup>1,4</sup>, Wenqi Huang<sup>1</sup>, Guohuan Song<sup>2,3</sup>, Guo-Huan Song<sup>1,4</sup>, Sijie Ma<sup>1</sup>,  
Zhiguang Cao<sup>5</sup>, Yue-Jiao Gong<sup>1\*</sup>

<sup>1</sup> School of Computer Science and Technology, South China University of Technology, Guangzhou, 510006, China;

<sup>2</sup> School of Computer Science and Technology, Zhejiang Normal University, Jinhua, 321004, China;

<sup>3</sup> Northern Computility, Beijing, China;

<sup>4</sup> Panorama Optimization, Guangzhou, China;

<sup>5</sup> School of Computing and Information Systems, Singapore Management University, 178902, Singapore;

\* Correspondence: [gongyuejiao@gmail.com](mailto:gongyuejiao@gmail.com)

## Abstract

Machine intelligence marks the ultimate dream of making machines' intelligence comparable to human beings. While recent progress in Large Language Models (LLMs) show substantial *specific skills* for a wide array of downstream tasks, they more or less fall shorts in *general intelligence*. Following correlation between intelligence and system 2 reasoning (slow thinking), in this paper, we aim to answering a worthwhile research question: could machine intelligence such as LLMs be evolved to acquire reasoning ability (not specific skill) just like our human beings? To this end, we propose evolutionary reasoning optimization (ERO) framework which performs *survival of the fittest* over a population of LLMs to search for individual with strong reasoning ability. Given a reasoning task, ERO first initializes multiple LLMs as a population, after which an evolutionary strategy evolves the population to maximize quantified reasoning score of the best individual. Based on experiments on representative test suites, we claim two surprising empirical discoveries: i) the latest LLMs such as GPT-5 still show limited system 2 reasoning ability; ii) with simple evolution-loop of ERO, a relatively weak model (Qwen-7B) could be enhanced to emerge powerful reasoning ability. Our project can be accessed at <https://github.com/MetaEvo/ERO> for reproduction needs.

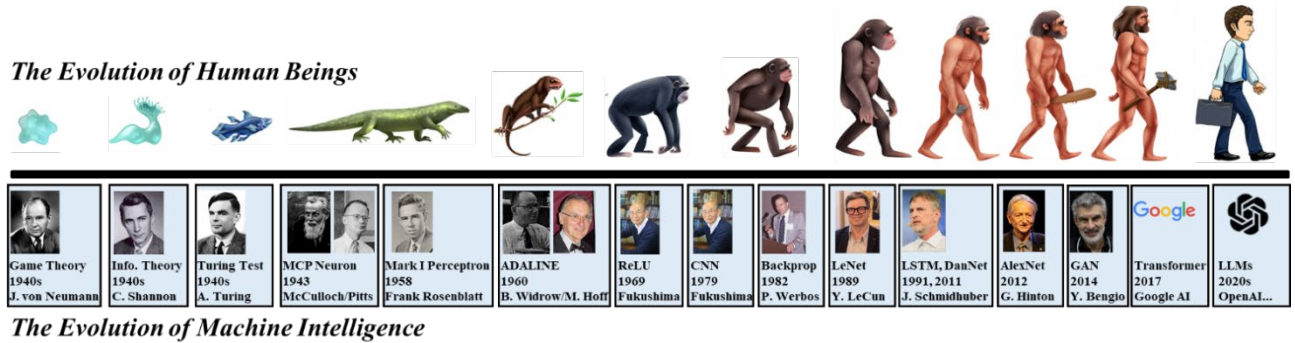
**Keywords:** Large Language Model; Logical Reasoning; System 2 Reasoning; Neuroevolution; Evolutionary Computation; Black-Box Optimization

---

## 1. Introduction

Machine intelligence (often interchangeably used with AI) has experienced ups and downs within a long river of history [1][2][3]. Since the initial proposal of AI at 1950s [4], an evolution path has been observed: from basic theories [5][6] to concrete architectures [7][8][9][10][11] and algorithms [12][13][14][15]. Today, the application of AI has spread to every corner of the world. Domains such as image processing [16], nature language processing [17] and scientific discovery [18] benefit from its learning power and corresponding human-competitive performance.

However, we should not overlook the dark side of advanced machine intelligence (i.e., LLMs) simply due to its twinkling academic and engineering achievements [19][20][21]. In other words, we have to realize that LLMs, though pre-trained with massive human knowledge prior, may still operate



**Figure 1:** A comparison between the evolution paths of human beings and machine intelligence.

at the pattern recognition (fast thinking, System 1 reasoning) level, and hence lacks long-chain, deep, logical reasoning ability (slow thinking, System 2 reasoning), as testified in recent competitions<sup>1</sup>.

As illustrated in Figure 1, such System 2 reasoning inability potentially roots from the essential difference between the evolution of machine intelligence and that of our human beings [22][23]. For human beings, we are continually involved in evolutionary process under open-ended environmental selection pressure, which follows the *survival of the fittest* principle proposed by Darwin [24]. The "open-ended" term is used to reference extreme generalization scenario where environmental uncertainty is naturally unknown by human beings [25]. In contrast, almost all machine intelligence instances are trained for specific application scopes explicitly restricted by their developers (human beings). The feedback or learning signal in their learning loops may inherently restrict them from general intelligence with logic reasoning [26]. To make this point clearer, we borrow the valuable perspective from developmental psychology [27], which holds the position that: human-level intelligence shows generalization and open-endedness and is capable of expanding far beyond its evolution path. More importantly, human is born with innate and evolution-driven knowledge priors such as elementary physics, goal-directness, arithmetic and geometry. These priors enable us to acquire certain skills efficiently [28], by System 2 slow thinking.

The gap between existing LLMs and general System 2 reasoning ability motivates us to explore possible solutions. An intuitive yet under-explored thought would be: Given that existing advanced LLMs have absorbed massive knowledge priors through pre-training with internet-scale corpus, can we further evolve them (e.g., Neuroevolution [29]) to attain System 2 reasoning ability? To answer this research question empirically, we in this paper propose **Evolutionary Reasoning Optimization (ERO)** framework that enables human-like evolution process for LLMs to adapt themselves in complex tasks that require System 2 reasoning. In our framework, the neural network parameters of a LLM are regarded as a holistic genotype space. At the beginning, given a complex reasoning task, a population of LLMs are randomly born via sampling from the genotype space. Then a  $(\mu + \lambda)$  Evolutionary Strategy (ES) [30][31] is applied to guide the LLM population toward more powerful System 2 reasoning performance on the target task. The evolution rule in our framework is purely objective-oriented: the LLM individual with higher reasoning ability survives and contributes to the reproduction of offspring, which is closely analogous to evolution of human beings. We provide an intuitive illustration in Figure 4 to showcase how ERO evolves a weak Qwen-7B model to surpass powerful GPT-5 model on reasoning tasks. We next provide a brief review of related works in Sec. 2, elaborate the technical details of **ERO** in Sec. 3 and discuss empirical results in Sec. 4 respectively.

<sup>1</sup> <https://arcprize.org/leaderboard>

## 2. Related Works

### 2.1. Reasoning in LLMs

Reasoning ability is regarded as a key for achieving human-level machine intelligence [32]. In particular, it relies on logical reasoning and systematic step-by-step thinking to ensure solving effectiveness on complex tasks, which is typically termed as System 2 reasoning. Compared to System 1 reasoning, which features fast, pattern recognition-based decision mapping, System 2 reasoning presents deliberate slow thinking, resulting in concise and rational problem solving via mitigating cognitive biases in System 1 reasoning. While the swift development of LLMs (e.g., DeepSeek-v3 [33], GPT-5 [34]) shows promising results on understanding and performing human-competitive tasks, they may still lack matched cognitive abilities with human beings in complex reasoning tasks [35].

To improve the capability of reasoning LLMs, initial exploration includes Chain-of-Thought (CoT) [36][37] and Tree-of-Thought (ToT) [38], which focus on preparing high-quality, step-by-step and fine-grained supervision data through decomposing the complex reasoning process into chain or tree structure. Given the data scaling difficulty and single-pass reasoning pattern in CoT and ToT, subsequent works further apply Monte Carlo Tree Search (MCTS) to allow LLMs revisit, reflect and refine their reasoning paths dynamically [39][40][41], or self-improvement strategies [42] that bootstrap training data from either iterative self-reflection [43] or rule-based reasoning path augmentation [44]. Beside these data curation designs, the training paradigm itself also plays crucial role in attaining robust reasoning LLMs. Common practice in up-to-date literature leans to reinforcement fine-tuning (RFT) with output reward modeling (ORM) [45] or process reward modeling (ORM) [46]. The former emphasizes scoring for final answer correctness and the latter pays efforts on fine-grained step-by-step reward labeling. Test time training (TTT) [47] is also adopted as effective post-training strategy to mitigate reasoning hallucination. For further reading, we suggest these surveys [32][48][49].

### 2.2. Reasoning LLMs Benchmarks

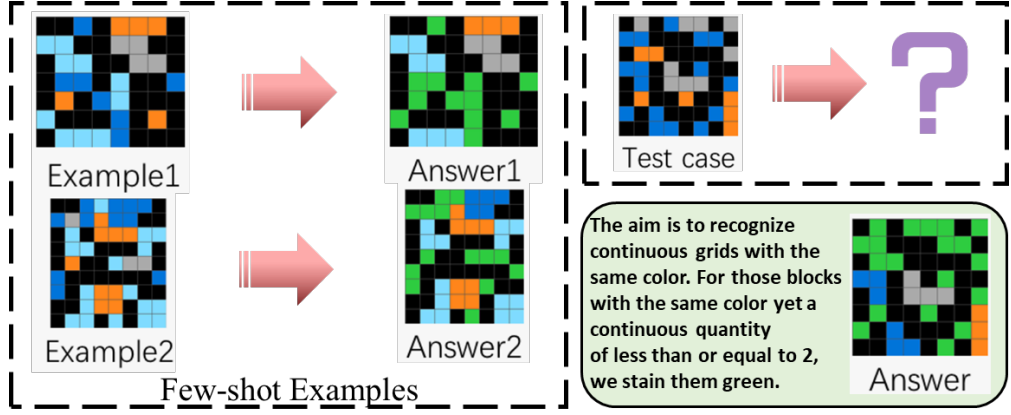
While recent advance of LLMs demonstrates that, with large-scale pre-training on massive and diverse corpus, these novel machine intelligence models rival or even surpass human's performance at specific testbeds [50][51][52], more evidences argue that they lack compositional System 2 reasoning ability in solving complex tasks as general intelligence [53]. To this end, a large body of related benchmarks have been curated to provide objective and challenging reasoning tasks for evaluating reasoning LLMs. According to their concrete task types, we could generally document them as: 1) Olympic-level mathematical reasoning benchmarks [54][55]; 2) Real world programming challenges summarized from GitHub [56]; 3) Scientific discovery process [57] in physics, chemistry, etc.; 4) Agentic automation workflow tests, e.g., constructing web application from zero [58]; 5) Human-level cognitive ability tests [28][35] that analog IQ examination.

In this paper, we focus on the last benchmark type, of which a representative benchmark is Abstraction and Reasoning Corpus (ARC) benchmark [28]. As illustrated in Figure 2, the testing task instance in ARC benchmark includes multiple few-shot examples and a test case for machine intelligence to solve, which stays close in format of psychometric intelligence test [59]. To figure out each puzzle, an intelligence must coherently enable its innate prior on object persistence and contact influence, goal-directedness, numbers and counting, etc., just like our human beings. According to the latest results, even the most powerful reasoning-reinforced LLMs (GPT 5 and Gemini 3) could only achieve scores no more than 55% on ARC-AGI-2 benchmark[35], with an evident reasoning gap against human panel (70%, according to [73]). The ARC benchmarks provide us a desirable testbed.

### 2.3. Evolutionary LLMs Enhancement

Evolutionary Algorithms (EAs) [60] are meta-heuristics that follow evolutionary principle in nature to optimize given problems through reproduction and selective pressure. Given EAs' high-level alignment with the evolution process of human beings and robust global optimization capability, they have been validated as powerful optimization techniques for many applications [61], except LLMs. Recently,

initial attempts have been made to explore the possibility of leveraging EAs to enhance LLMs' performances. While limited, these efforts have seen delightful effects such as prompt optimization through textual evolution [62], program evolution through LLM-level genetic programming [21][63], novel ability composition through model merge recipes [64][65], incremental and dynamic prompting through evolutionary context engineering [66]. However, to the best of our knowledge, none of prior works focus on the core vision of LLMs: reasoning like human beings. This highlights the motivation of our paper.



**Figure 2:** A reasoning task example in ARC benchmark.

**Table 1:** Pass@1 scores of LLMs baselines across 15 ARC tasks, with their task properties attached at the top of the table.

Tasks	T1	T2	T3	T4	T5	T6	T7	T8	T9	T10	T11	T12	T13	T14	T15
<b>Properties</b>															
Object cohesion	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✗	✗	✓
Object persistence	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✓	✗
Object influence via contact	✗	✗	✗	✓	✗	✓	✗	✗	✓	✗	✗	✓	✗	✗	✗
Goal-directedness	✓	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗
Numbers and Counting	✗	✗	✗	✗	✗	✓	✓	✓	✗	✗	✓	✓	✗	✗	✓
Basic Geometry and Topology	✓	✓	✓	✓	✗	✓	✗	✗	✓	✓	✗	✓	✓	✓	✓
<b>Performance</b>															
Ours (ERO+Qwen2.5-7B)	0.7765	<b>0.7820</b>	<b>0.9845</b>	0.7828	<b>0.9016</b>	<b>1.0000</b>	<b>0.8100</b>	<b>1.0000</b>	0.9461	<b>0.8182</b>	<b>1.0000</b>	0.9627	0.7193	<b>0.7315</b>	0.7073
Qwen2.5-7B	0.7059	0.1132	0.9380	0.4253	0.2623	0.9344	0.3584	0.6400	0.8491	0.3333	0.6400	0.9379	0.6486	0.5370	0.6098
Qwen2.5-32B	0.6118	0.3831	0.9535	0.4143	0.5164	0.6393	0.0924	<b>1.0000</b>	0.4315	0.4848	0.5200	0.4752	<b>0.7235</b>	0.6205	0.4268
GPT-4o-mini	0.1401	0.1200	0.9380	0.3875	0.3115	0.9344	0.3122	0.6400	0.9212	0.4545	0.6400	0.9379	0.5489	0.5507	0.6341
GPT-4o	0.6195	0.2699	<b>0.9767</b>	0.4434	0.2295	0.7541	<b>0.6380</b>	<b>1.0000</b>	0.8880	<b>0.6061</b>	0.6400	<b>0.9689</b>	0.6861	<b>0.7205</b>	0.6341
GPT-5	<b>0.8647</b>	<b>0.6505</b>	0.4961	<b>1.0000</b>	<b>0.8934</b>	<b>1.0000</b>	0.4027	<b>1.0000</b>	<b>0.9647</b>	0.5455	<b>0.8200</b>	0.9472	0.4376	0.3260	<b>0.7195</b>

### 3. Evolutionary Reasoning Optimization

In this section, we elaborate both the general picture and specific designs of our ERO framework to clarify how we address reasoning enhancement for LLMs via EAs perspective. Generally speaking, ERO operates as a neuroevolution [29] approach, which is under the umbrella of evolutionary strategy (ES) [30][31] framework. We present the overall workflow of ERO in Alg. 1, where starting from an existing LLM, an iterative searching process is deployed to evolve the parameters of the LLM toward high reasoning performance on the given reasoning task. However, we must note that it is neither practical nor efficient to run ERO in a standard ES procedure. We next detail the key challenges and corresponding tailored designs in ERO.

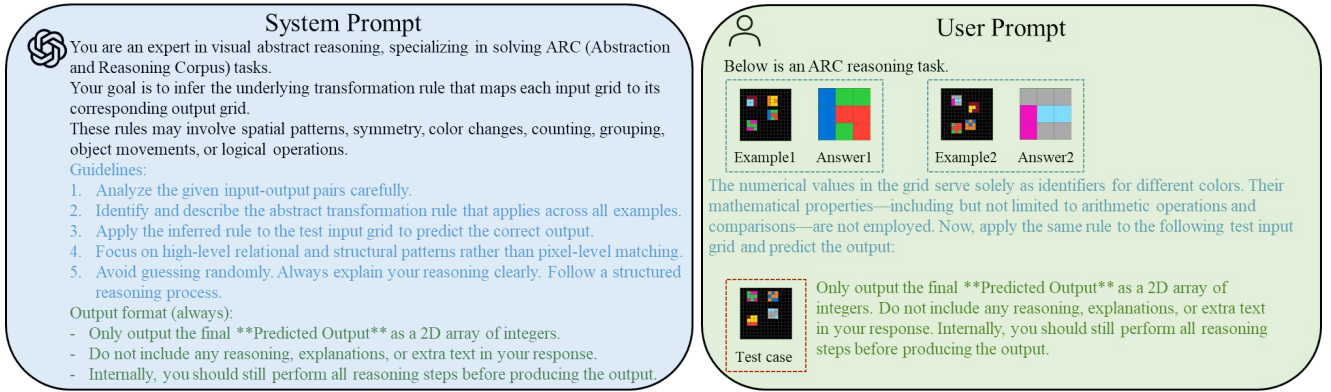
**Sampling Strategy:** In ERO, we have to first determine sampling strategy (i.e., mean and covariance parameters) to serve as initialization module and hence kick out subsequent evolution process. For the mean parameter  $\mu$ , we could simply set it as the weights of the LLM (denoted as  $\theta^0$ ). The real challenge is how to determine covariance parameter  $\Sigma$ . Although ES has been previously used in evolving relatively smaller neural networks [29][67], where the value ranges of parameters are



controllable and hence we could set identical entries for  $\Sigma$  matrix, it is absolutely not the case in LLMs. This is backed up by our preliminary experiment, where we conducted a statistical summary on different LLMs and found out that the value ranges of LLMs' parameters vary a lot. However, we also found out that the value ranges of layer-wise parameters are more stable. Based on such observation, we determine the entries of variance matrix  $\Sigma$  by the principle below:

$$\Sigma[k, k] = \epsilon \times \frac{1}{|L_k|} \sum_{n=1}^{|L_k|} \theta^{(0)}[L_k][n] \quad (1)$$

where  $k$  denotes  $k$ -th neural network parameters of the selected LLM,  $L_k$  is the network layer where  $k$ -th parameter locate at,  $\epsilon$  is value between 0~1 to control the variance strength,  $[\cdot]$  is the indexing operation. We set  $\Sigma$  once leave it fixed until the end. A population of LLMs with the same architecture with  $\theta^{(0)}$  are then sampled by the constructed gaussian distribution (line 4 in Alg. 1).



**Figure 3:** System prompt and User prompt we used across all baselines.

**Scoring Function:** Given a population of  $\lambda$  sampled LLMs at  $g$ -th generation:  $\{\theta^{(g),i}\}_{i=1}^{\lambda}$ , the underlying ES process in ERO needs proper evaluation metric (scoring function) to measure the reasoning performances of these LLM individuals on the given task  $\tau$ . A general form of such scoring function can be formulated as  $\mathbb{S}(\theta|\tau)$ , where  $\theta$  denotes a tested LLM. We would like to clarify that our ERO does not restrict concrete implementation of the scoring function, instead, it can be quite flexible to tailor appropriate scoring schemes for different reasoning tasks. One can surely use generic schemes such as process reward model [46] that regards any reasoning task as standard reasoning chain and credits those matched reasoning steps. On the other hand, one can also customize special  $\mathbb{S}$  function for specific task. Since our ERO is a purely objective-oriented optimization system, all it need is a scalar objective to minimize or maximize. We take the testbed we select for this paper (ARC benchmark) as an example. In ARC, the answer of a reasoning task instance is typically a 1-D or 2-D array indicating the colors of grids. By representing them as strings, one can simply compute the score as:

$$\mathbb{S}(\theta|\tau) = 1 - \frac{\text{lev}(\hat{A}(\tau|\theta), A(\tau))}{\max(\text{Len}(A(\tau)), \text{Len}(\hat{A}(\tau|\theta)))} \quad (2)$$

where  $\text{lev}(\cdot, \cdot)$  is the Levenshtein distance [68] between two strings,  $A$  and  $\hat{A}$  is the ground truth and predicted answer respectively. In this paper, ERO aims at maximizing the LLM's performance on ARC tasks through maximizing corresponding scoring function values.

**Island Architecture:** Given the massive searching space of LLM's neural network parameters, island-based population architecture could be a useful strategy to enhance the searching diversity of underlying ES process, which may further improve the final optimization performance [69]. Besides, since LLMs are inherently aligned with multi-card computational resources and distributed

computational methods, island architecture is a coherent choice in LLM-based evolution frameworks [70][71]. To this end, our ERO instantiates multiple LLM populations as independent islands, which sample and evaluate LLM individuals (lines 4~5 of Alg. 1) in parallel. The communication (fitness aggregation) across different islands occurs when we have to aggregate elite LLM individuals and accordingly update the mean and variance parameters of ES process (lines 6~7 of Alg. 1). Unlike vanilla ES, the  $\mu$  elite LLM individuals are selected as the  $\lfloor \frac{\mu}{Z} \rfloor$  best individuals per island, where  $Z$  is the number of islands deployed. Once the elite individuals are voted out, we update the mean parameters used for next-generation sampling by averaged aggregation. Note that we keep a fixed variance matrix  $\Sigma$  to maintain continuous exploration strength along the evolution process.

**Ray Acceleration:** As we mentioned above, the island architecture allows us to incorporate advanced distributed ML techniques to reduce the running complexity of LLM-based evolution frameworks. This is particularly useful in our ERO, since the scoring evaluation is actually time-consuming, where each LLM individual is fed with reasoning questions and prompted to output reasoning steps and answers. We hence introduce Ray<sup>2</sup>, a large-scale ML-enabled computational framework, to distribute each island in ERO onto a separate GPU of a multi-GPU computer/cluster. The Ray parallelism not only enables distributed island-based evolution, but also further facilitates fine-grained parallel evaluation within each island, reducing the running time of ERO from days to hours. In practice, the concrete parallel degree varies due to different hardware conditions.

---

**Algorithm 1:** Evolutionary Reasoning Optimization

---

Input: LLM  $\theta^0$ ; reasoning task  $\tau$ ; population size  $\lambda$ ; elite group size  $\mu$ ; optimization budget  $G$ .

Output: best LLM individual  $\theta^*$  found ever.

- 1: Attain layer-wise covariance  $\Sigma$  from  $\theta^0$
  - 2: Let  $g = 1$
  - 3: **while**  $g < G$  **do**
  - 4:   Sample  $\lambda$  LLMs:  $\{\theta^{(g),i}\}_{i=1}^{\lambda} \sim N(\theta^{(g-1)}, \Sigma)$
  - 5:   Evaluate their reasoning scores:  $\{\mathbb{S}(\theta^{(g),i} | \tau)\}_{i=1}^{\lambda}$
  - 6:   Select  $\mu$  top-scoring LLMs:  $\{\hat{\theta}^{(g),j}\}_{j=1}^{\mu}$
  - 7:   Update  $\theta^{(g)} = \frac{1}{\mu} \sum_{j=1}^{\mu} \hat{\theta}^{(g),j}$
  - 8:    $g = g + 1$
  - 9: **end while**
  - 10: **return** the LLM individual with the best score
- 

**Cache Optimization:** One may question about how could a large population of LLMs be loaded within a single 4-GPU or 8-GPU computer/cluster, since a single LLM may require at least 10~20 GB GPU memory. The solution we propose is to subtly and flexibly leverage limited cache memory. In specific, we only maintain necessary LLM information in an *on-the-fly* fashion for each island (i.e., each GPU node). The necessary LLM information includes the mean parameters at current optimization generation ( $\theta^{(g)}$ ), the layer-wise variance matrix  $\Sigma$ , the elite pool used for maintaining  $\lfloor \frac{\mu}{Z} \rfloor$  elite LLM individuals. The elite pool is dynamically updated when a newly sampled LLM individual gets better reasoning score than those in the pool, where the older elite is replaced by the newly sampled one. With such cache memory optimization, ERO could evolve hundreds of LLM individuals simultaneously on GPU memory-limited platform.

---

<sup>2</sup> <https://github.com/ray-project/ray>

## 4. Empirical Validation and Discussion

### 4.1. Experimental Setup

We list detailed settings of each part in ERO here, which could be generally divided into three categories:

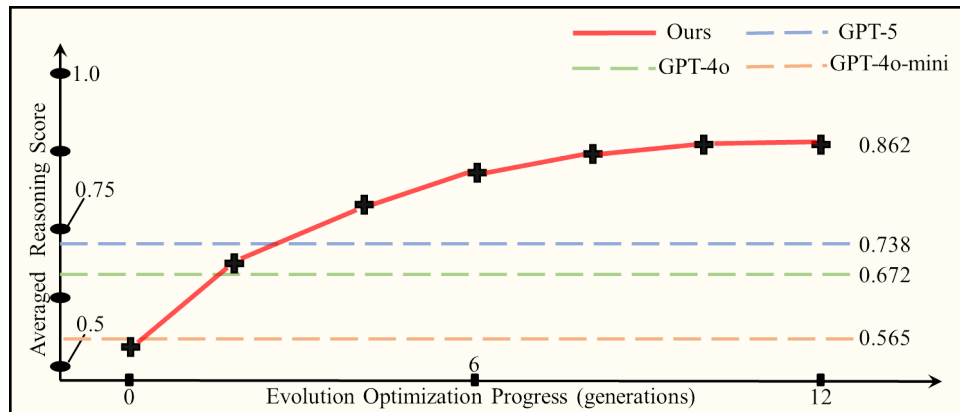
**ERO's Settings:** We select *Qwen-7B*<sup>3</sup> as the initial LLM  $\theta^{(0)}$  to be evolved. The reason behind such selection is that this relatively poor-reasoning model could facilitate validation on effectiveness of our ERO. For the hyper-parameters of the underlying island-based ES, we set its population size  $\lambda = 1000$  which are evenly distributed to  $Z = 4$  islands, elite pool size  $\mu = 4$  and the optimization budget  $G = 12$  generations. All experiments are run on a high-performance instance of a GPU cluster, which comprises an Intel Xeon 8558P CPU, 128 GB RAM and 4×64 GB virtual GPU nodes based on Nvidia H20 GPU.

**Testbed:** As a preliminary study and due to limited computational resources, in this paper, we have randomly sampled 15 reasoning task instances from hundreds of instances in ARC-1 benchmark [28]. We mark these 15 tested instances as  $T1 \sim T15$ . We present at upper half of Table 1 the fine-grained properties of these instances in terms of their correspondence to innate cognitive abilities of human beings. Refer to our project for their correspondence to ARC-1 indices and concrete task descriptions and visualizations.

**Baselines:** We include 6 baselines in the comparison experiments: 1) *Ours*: the *Qwen-7B* model evolved by our ERO on the given ARC-1 reasoning task instance; 2) *Qwen-7B*: the same *Qwen-7B* pre-trained checkpoint, without ERO's evolution; 3) *Qwen-32B*<sup>4</sup>: a much larger *Qwen* model with stronger general task solving ability than the 7B model; 4) *GPT-4o-mini*<sup>5</sup>, 5) *GPT-4o*<sup>6</sup> and 6) *GPT-5*<sup>7</sup>, which are three GPT-series models enhanced with multi-modal processing ability and chain-based reasoning capability. For *Ours* and *Qwen-7B*, we deploy their checkpoints at our local GPU server. For the rest of baselines, we call their corresponding APIs. Their key hyper-parameters such as temperature and top-p sampling rate follow default values. For GPT-5, we use its default reasoning efforts level ("minimal").

### 4.2. Major Results

For all of the selected baselines, we use a pre-designed standard prompt template to ensure fair evaluation, as illustrated in Figure 3. By using this standard template, we could test selected baselines on the 15 reasoning tasks sampled from ARC-1 benchmark, and then compute their per-instance pass@1 reasoning scores (as we defined in Eq. (2)). We next present these results and corresponding discussions.



**Figure 4:** Evolution curve of ERO on ARC benchmark.

<sup>3</sup> <https://huggingface.co/Qwen/Qwen2.5-VL-7B-Instruct>

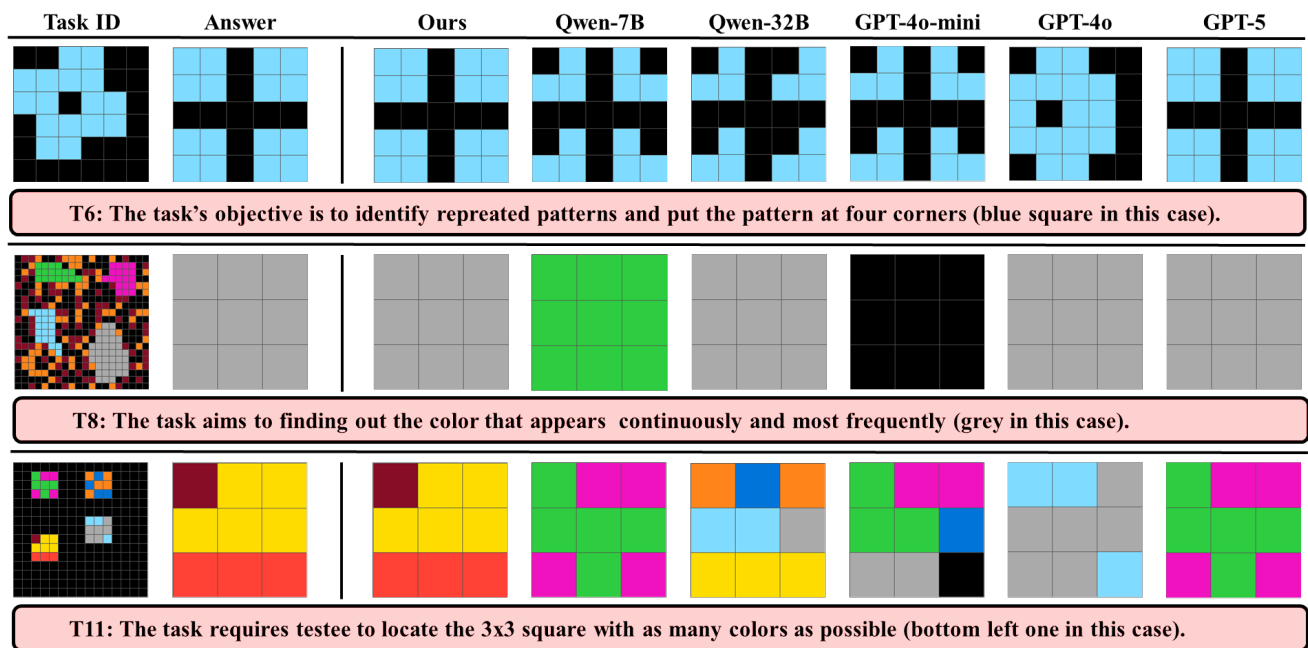
<sup>4</sup> <https://dashscope.aliyuncs.com/compatible-mode/v1>

<sup>5</sup> <https://platform.openai.com/docs/models/gpt-4o-mini>

<sup>6</sup> <https://platform.openai.com/docs/models/gpt-4o>

<sup>7</sup> <https://platform.openai.com/docs/models/gpt-5>

**Evolutionary Convergence:** We first demonstrate the effectiveness of our ERO by illustrating its evolutionary convergence curve as shown in Figure 4, where each scalar point in the red line is the average reasoning score of ERO across 15 tested reasoning tasks. We also attach the average scores of three advanced GPT-series baselines with dashed lines. The results in Figure 4 demonstrate that while a simple pre-trained Qwen-7B model underperforms the GPT models due to its limited capacity and pre-training data scale, it could be evolved by our ERO to surpass these advanced baselines on reasoning tasks. This finding may also indicate the knowledge prior redundancy of existing LLMs. We may not need continually scale both the model capacity and training data size to enable LLM's human-level reasoning ability. On the contrary, such ability may conceal itself within the LLM's parameters, and could be adapted to specific reasoning task through evolution. The results above at least demonstrate potential of evolutionary algorithms on LLM's post-tuning.



**Figure 5:** Showcases on the effectiveness our ERO for boosting the understanding and reasoning ability of LLMs.

**Performance Comparison:** We further present the per-instance performance comparison between our ERO and other baselines in the lower half of Table 1. Where the best and second-best are labeled in bold and underlined respectively. We also specifically mark the results of our ERO and Qwen-7B in light blue to highlight the relative improvement. From the results, we can observe that: 1) ERO could significantly improve the reasoning capability of Qwen-7B through 12 evolution generations, which cross-validates that intelligence (whether organic or machine-based) obeys evolution principle (*survival-of-the-fittest*); 2) With our ERO, a relatively weak Qwen-7B LLM could be evolved to perform competitively with one of the most advanced LLMs: GPT-5. On 8 of the 15 tested task instances, ERO presents significant performance advantage; 3) The reasoning capability of LLMs may not root from existing scaling law in training these LLMs. Direct evidence lies in the comparison between Qwen-7B and Qwen-32B models. On 8 of the 15 reasoning tasks, a smaller Qwen-7B model presents better logical reasoning and understanding level than its "improved version". This might indicate that we should pay more attention on multi-dimensional solutions for reasoning enhancement of LLMs, not only the scale of LLM pre-training.

We also showcase in Figure 5 three task instances (T6, T8 and T11) where our ERO successfully evolves the initial Qwen-7B model from completely wrong reasoning to crystal correct answer. As their



descriptions and visualizations presented in the figure, these ARC-1 reasoning tasks challenge the innate abilities of intelligence of human beings, let alone the LLMs never being trained on such tasks.

### 4.3. An Important Future Work

In this paper, we mainly focus on the evolution of LLM's reasoning capability under a given reasoning task. While the results mentioned in previous sections have clearly demonstrated that introducing evolutionary perspective into LLM's intelligence enhancement could result in surprising and promising effects, we have to note that the evolution of human beings may not be such simple, i.e., in an *adaption-per-task* fashion. On the contrary, the subtle evolution of human beings emerges in the remix of complex environmental dynamics and concurrent multitasking. This outlines an important and promising future work of our ERO, which is the meta-evolution across reasoning task distribution:

$$S_{meta} = \mathbb{E}_{\tau \sim \Omega}[S(\theta|\tau)]$$

which is the expectation of reasoning scores over a reasoning task distribution  $\Omega$ . As computing power and evolution paradigm (e.g., [72]) continue to iterate and update, we may witness in the near future the emergence of machine intelligence species with diverse behavior and characteristics (e.g., "The Matrix" movie), purely by evolution.

## 5. Conclusion

The position of this paper bridges the evolutionary computation community and LLMs community by proposing the ERO framework, which iteratively evolves LLM's parameters to maximize its System 2 reasoning scores on given reasoning tasks. At the core of ERO, we introduce island architecture-based evolutionary strategy to ensure searching diversity and quality, which attains reasoning performance gain effectively. Combined with specially designed cache optimization and ray acceleration mechanisms, ERO is capable of evolving a large population of LLMs on relatively limited computational resources. We validate ERO's potential by comparing it to existing representative LLMs on ARC benchmark. The promising results not only demonstrate evolution of LLMs is useful for intelligence enhancement, but may also reveal implicit connections between organic human beings and connectionism-based machine intelligence. We hope this work could appeal for more research efforts on evolutionary machine intelligence, and more importantly, exploration on more possibilities.

## 6. Acknowledgments

We acknowledge that the source material we used in Figure 1 are partly designed by Freepic and Wikipedia. We also appreciate the editors in JISS for their timely review and step-by-step instructions to help us improve our paper.

### Author Contribution

**Conceptualization** (Zeyuan Ma), **Data curation** (Wenqi Huang), **Formal analysis** (Yuejiao Gong), **Finding acquisition** (Zeyuan Ma), **Investigation** (Zeyuan Ma, Wenqi Huang), **Methodology** (Zeyuan Ma), **Project administration** (Wenqi Huang), **Resources** (Guohuan Song, Zhiguang Cao), **Software** (Wenqi Huang), **Supervision** (Yuejiao Gong, Zhiguang Cao), **Validation** (Hongshu Guo, Sijie Ma), **Visualization** (Wenqi Huang), **Writing – original draft** (Zeyuan Ma), **Writing – review & editing** (Yuejiao Gong, Hongshu Guo, Sijie Ma).

### Funding

This research was supported by National Natural Science Foundation of China (Grant No. 62276100).

## References

1. Legg, Shane, and Marcus Hutter. "Universal intelligence: A definition of machine intelligence." *Minds and machines* 17.4 (2007): 391-444.
2. Minsky, Marvin. "Steps toward artificial intelligence." *Proceedings of the IRE* 49.1 (2007): 8-30.
3. LeCun, Yann, Yoshua Bengio, and Geoffrey Hinton. "Deep learning." *nature* 521.7553 (2015): 436-444.
4. McCarthy, John, et al. "A proposal for the Dartmouth summer research project on artificial intelligence, august 31, 1955." *AI magazine* 27.4 (2006): 12-12.
5. Shannon, Claude E. "A mathematical theory of communication." *The Bell system technical journal* 27.3 (1948): 379-423.
6. Turing, Alan M. "Computing machinery and intelligence (1950)." *Mind* 59.236 (2021): 33-60.
7. Rosenblatt, Frank. "The perceptron: a probabilistic model for information storage and organization in the brain." *Psychological review* 65.6 (1958): 386.
8. Fukushima, Kuniyoshi. "Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position." *Biological cybernetics* 36.4 (1980): 193-202.
9. Hochreiter, Sepp, and Jürgen Schmidhuber. "Long short-term memory." *Neural computation* 9.8 (1997): 1735-1780.
10. Vaswani, Ashish, et al. "Attention is all you need." *Advances in neural information processing systems* 30 (2017).
11. Gu, Albert, and Tri Dao. "Mamba: Linear-time sequence modeling with selective state spaces." *First conference on language modeling*. 2024.
12. Robbins, Herbert, and Sutton Monro. "A stochastic approximation method." *The annals of mathematical statistics* (1951): 400-407.
13. Werbos, Paul John. *The roots of backpropagation: from ordered derivatives to neural networks and political forecasting*. John Wiley & Sons, 1994.
14. Loshchilov, Ilya, and Frank Hutter. "Decoupled weight decay regularization." *arXiv preprint arXiv:1711.05101* (2017).
15. Graves, Alex. "Generating sequences with recurrent neural networks." *arXiv preprint arXiv:1308.0850* (2013).
16. Gonzalez, Rafael C. *Digital Image Processing*. Pearson Education India, 2009.
17. Bengio, Yoshua, et al. "A neural probabilistic language model." *Journal of machine learning research* 3.Feb (2003): 1137-1155.
18. Jumper, John, et al. "Highly accurate protein structure prediction with AlphaFold." *nature* 596.7873 (2021): 583-589.
19. Zhou, Hao, et al. "Large language model (LLM) for telecommunications: A comprehensive survey on principles, key techniques, and opportunities." *IEEE Communications Surveys & Tutorials* 27.3 (2024): 1955-2005.
20. Li, Yuanchun, et al. "Personal LLM agents: Insights and survey about the capability, efficiency and security." *arXiv preprint arXiv:2401.05459* (2024).
21. Novikov, Alexander, et al. "AlphaEvolve: A coding agent for scientific and algorithmic discovery." *arXiv preprint arXiv:2506.13131* (2025).
22. Pinker, Steven. *The blank slate: The modern denial of human nature*. Penguin, 2003.
23. Cosmides, Leda, and John Tooby. "Origins of domain specificity: The evolution of functional organization." *Mapping the mind: Domain specificity in cognition and culture* 853116 (1994).
24. Darwin, Charles, John Wyon Burrow, and John Wyon Burrow. *The origin of species by means of natural selection: or, the preservation of favored races in the struggle for life*. New York: AL Burt, 2009.
25. Wolpert, David H. "What can we know about that which we cannot even imagine?." *New Frontiers in Science in the Era of AI*. Cham: Springer Nature Switzerland, 2024. 301-331.
26. Wolpert, David H., and William G. Macready. "No free lunch theorems for optimization." *IEEE transactions on evolutionary computation* 1.1 (2002): 67-82.

27. Spelke, Elizabeth S., and Katherine D. Kinzler. "Core knowledge." *Developmental science* 10.1 (2007): 89-96.
28. Chollet, François. "On the measure of intelligence." *arXiv preprint arXiv:1911.01547* (2019).
29. Stanley, Kenneth O., et al. "Designing neural networks through neuroevolution." *Nature Machine Intelligence* 1.1 (2019): 24-35.
30. Rechenberg, Ingo. "Evolutionsstrategien." *Simulationsmethoden in der Medizin und Biologie: Workshop, Hannover, 29. Sept.–1. Okt. 1977*. Berlin, Heidelberg: Springer Berlin Heidelberg, 1978.
31. Beyer, Hans-Georg, and Hans-Paul Schwefel. "Evolution strategies—a comprehensive introduction." *Natural computing* 1.1 (2002): 3-52.
32. Li, Zhong-Zhi, et al. "From system 1 to system 2: A survey of reasoning large language models." *arXiv preprint arXiv:2502.17419* (2025).
33. Liu, Aixin, et al. "Deepseek-v3 technical report." *arXiv preprint arXiv:2412.19437* (2024).
34. OpenAI. "GPT-5 System Card." 2025, <https://cdn.openai.com/gpt-5-system-card.pdf>
35. Chollet, Francois, et al. "Arc-agi-2: A new challenge for frontier ai reasoning systems." *arXiv preprint arXiv:2505.11831* (2025).
36. Wei, Jason, et al. "Chain-of-thought prompting elicits reasoning in large language models." *Advances in neural information processing systems* 35 (2022): 24824-24837.
37. Kim, Seungone, et al. "The cot collection: Improving zero-shot and few-shot learning of language models via chain-of-thought fine-tuning." *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*. 2023.
38. Yao, Shunyu, et al. "Tree of thoughts: Deliberate problem solving with large language models." *Advances in neural information processing systems* 36 (2023): 11809-11822.
39. Li, Qingyao, et al. "Rethinkmcts: Refining erroneous thoughts in Monte Carlo tree search for code generation." *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*. 2025.
40. Cheng, Jiale, et al. "Spar: Self-play with tree-search refinement to improve instruction-following in large language models." *arXiv preprint arXiv:2412.11605* (2024).
41. Zhao, Yu, et al. "Marco-o1: Towards open reasoning models for open-ended solutions." *arXiv preprint arXiv:2411.14405* (2024).
42. Huang, Jiaxin, et al. "Large language models can self-improve." *Proceedings of the 2023 conference on empirical methods in natural language processing*. 2023.
43. Zelikman, Eric, et al. "Star: Bootstrapping reasoning with reasoning." *Advances in Neural Information Processing Systems* 35 (2022): 15476-15488.
44. Guan, Xinyu, et al. "rStar-Math: Small LLMs Can Master Math Reasoning with Self-Evolved Deep Thinking." *arXiv preprint arXiv:2501.04519* (2025).
45. Cobbe, Karl, et al. "Training verifiers to solve math word problems." *arXiv preprint arXiv:2110.14168* (2021).
46. Lightman, Hunter, et al. "Let's verify step by step." *The Twelfth International Conference on Learning Representations*. 2023.
47. Yang, Ling, et al. "Reasonflux: Hierarchical LLM reasoning via scaling thought templates." *arXiv preprint arXiv:2502.06772* (2025).
48. Huang, Jie, and Kevin Chen-Chuan Chang. "Towards reasoning in large language models: A survey." *Findings of the association for computational linguistics: ACL 2023*. 2023.
49. Chen, Qiguang, et al. "Towards reasoning era: A survey of long chain-of-thought for reasoning large language models." *arXiv preprint arXiv:2503.09567* (2025).
50. Wang, Alex, et al. "Superglue: A stickier benchmark for general-purpose language understanding systems." *Advances in neural information processing systems* 32 (2019).
51. Adiwardana, Daniel, et al. "Towards a human-like open-domain chatbot." *arXiv preprint arXiv:2001.09977* (2020).

52. Huang, Beichen, Ran Cheng, and Kay Chen Tan. "EvoGit: Decentralized Code Evolution via Git-Based Multi-Agent Collaboration." *arXiv preprint arXiv:2506.02049* (2025).
53. Lee, Seungpil, et al. "Reasoning abilities of large language models: In-depth analysis on the abstraction and reasoning corpus." *ACM Transactions on Intelligent Systems and Technology* (2024).
54. He, Chaoqun, et al. "Olympiadbench: A challenging benchmark for promoting AGI with olympiad-level bilingual multimodal scientific problems." *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 2024.
55. Yao, Jian, Ran Cheng, and Kay Chen Tan. "VAR-MATH: Probing True Mathematical Reasoning in LLMs via Symbolic Multi-Instance Benchmarks." *arXiv preprint arXiv:2507.12885* (2025).
56. Jimenez, Carlos E., et al. "Swe-bench: Can language models resolve real-world GitHub issues?." *arXiv preprint arXiv:2310.06770* (2023).
57. Wang, Yubo, et al. "MMLU-pro: A more robust and challenging multi-task language understanding benchmark." *Advances in Neural Information Processing Systems* 37 (2024): 95266-95290.
58. Zhou, Shuyan, et al. "Webarena: A realistic web environment for building autonomous agents." *arXiv preprint arXiv:2307.13854* (2023).
59. Carpenter, Patricia A., Marcel A. Just, and Peter Shell. "What one intelligence test measures: a theoretical account of the processing in the Raven Progressive Matrices Test." *Psychological review* 97.3 (1990): 404.
60. Jin, Yaochu, and Jürgen Branke. "Evolutionary optimization in uncertain environments-a survey." *IEEE Transactions on evolutionary computation* 9.3 (2005): 303-317.
61. Slowik, Adam, and Halina Kwasnicka. "Evolutionary algorithms and their applications to engineering problems." *Neural Computing and Applications* 32.16 (2020): 12363-12379.
62. Guo, Qingyan, et al. "Connecting large language models with evolutionary algorithms yields powerful prompt optimizers." *arXiv preprint arXiv:2309.08532* (2023).
63. Liu, Fei, et al. "Evolution of heuristics: Towards efficient automatic algorithm design using large language model." *arXiv preprint arXiv:2401.02051* (2024).
64. Akiba, Takuya, et al. "Evolutionary optimization of model merging recipes." *Nature Machine Intelligence* 7.2 (2025): 195-204.
65. Abrantes, João, Robert Lange, and Yujin Tang. "Competition and Attraction Improve Model Fusion." *Proceedings of the Genetic and Evolutionary Computation Conference*. 2025.
66. Zhang, Qizheng, et al. "Agentic context engineering: Evolving contexts for self-improving language models." *arXiv preprint arXiv:2510.04618* (2025).
67. Wierstra, Daan, et al. "Natural evolution strategies." *The Journal of Machine Learning Research* 15.1 (2014): 949-980.
68. Lcvenshtcin, V. I. "Binary coors capable or 'correcting deletions, insertions, and reversals." *Soviet physics-doklady*. Vol. 10. No. 8. 1966.
69. Gong, Yue-Jiao, et al. "Distributed evolutionary algorithms and their models: A survey of the state-of-the-art." *Applied Soft Computing* 34 (2015): 286-300.
70. Romera-Paredes, Bernardino, et al. "Mathematical discoveries from program search with large language models." *Nature* 625.7995 (2024): 468-475.
71. Lee, Kuang-Huei, et al. "Evolving deeper LLM thinking." *arXiv preprint arXiv:2501.09891* (2025).
72. Sarkar, Bidipta, et al. "Evolution Strategies at the Hyperscale." *arXiv preprint arXiv:2511.16652* (2025).
73. LeGris, Solim, et al. "H-ARC: A robust estimate of human performance on the abstraction and reasoning corpus benchmark." *arXiv preprint arXiv:2409.01374* (2024).