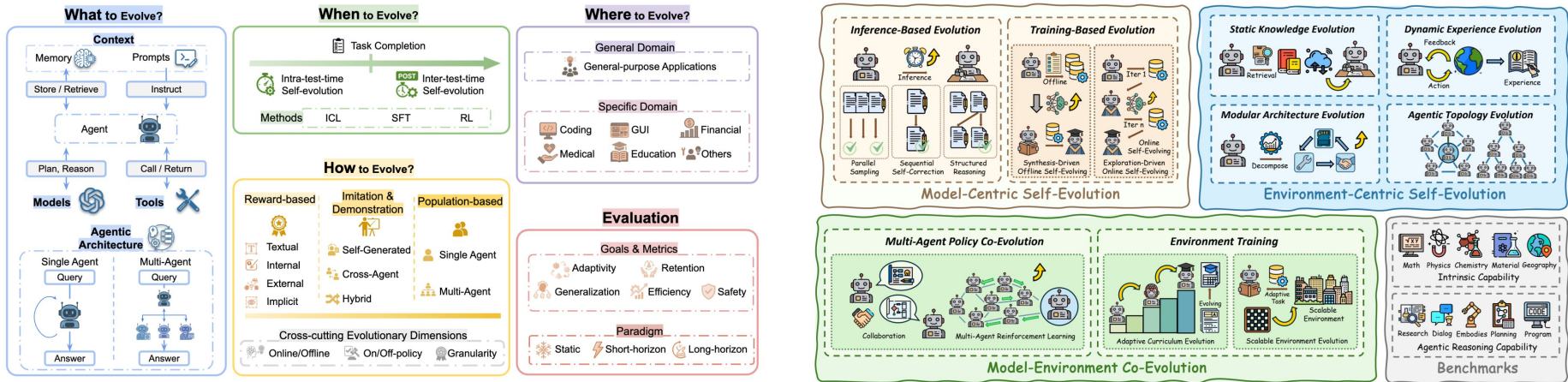


Position: Self-Play Only Evolves When Self-Synthetic Pipeline Ensures Learnable Information Gain

Revisit Self-Play and Self-Evolve From A Computational Complexity and Information-Theory Perspective

Wei Liu, Siya Qi, Yali Du, Yulan He

Self-Evolution



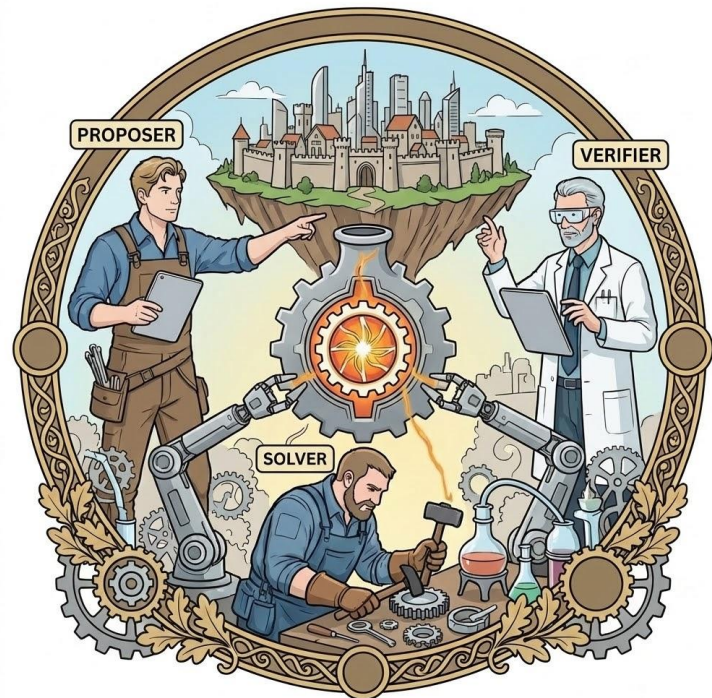
Various Definition:

- Continual Learning
- **Zero-data Self-Training**
- Prompting (update memory and experience)

Triadic Roles for Self-Evolution on Models

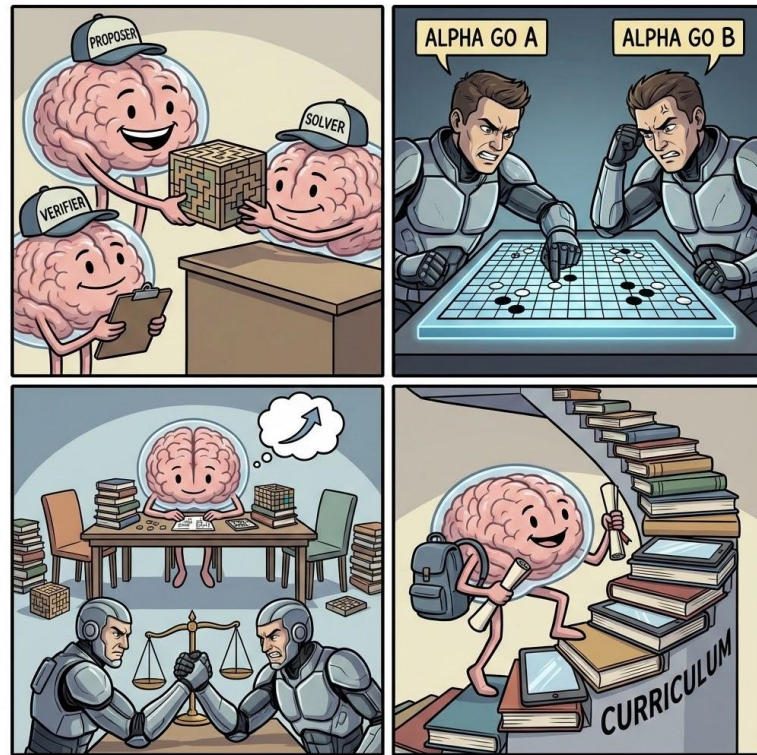
- **Proposer**, who propose the task
 - **Solver**, who gives the solution for the task
 - **Verifier**, who verify the solution
-
- Triadic for one (**P+S+V**) means all these three roles are put in one LLM
 - So co(tri)-evolution becomes self-evolution

$$\theta^* = \arg \max_{\theta} V_{\theta}(P_{\theta}, S_{\theta}(P_{\theta}))$$



Not Real Self-Play

- Not zero-sum game
- No Nash Equilibrium
- More like autonomous Curriculum Learning
- So, what do we expect from such LLM self-play?



It's self-synthetic data stream!

- Model *synthesis* its training data
- Model *synthesis* its feedback signal
- Model is trained on all these *synthetic* stuff

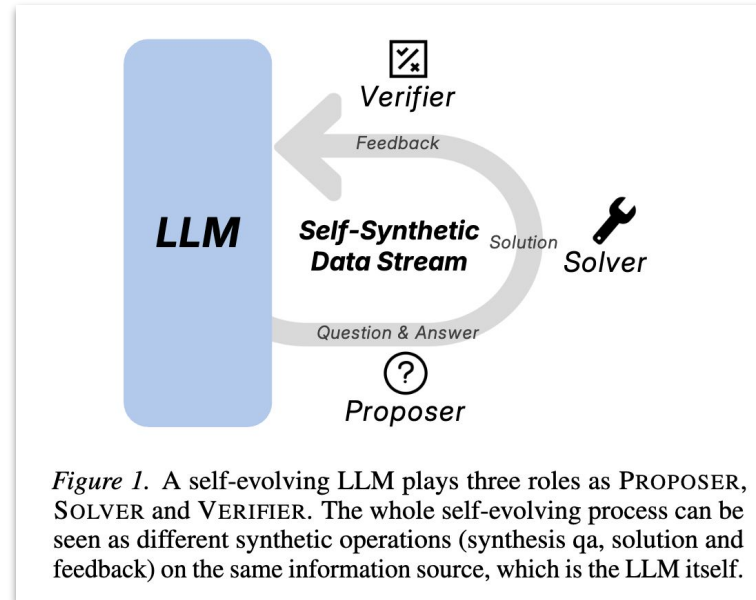
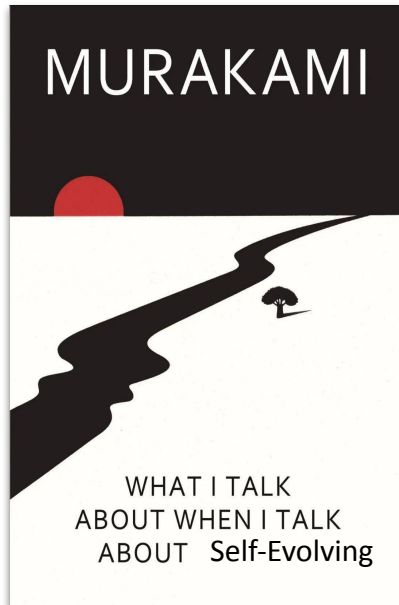


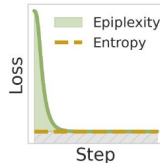
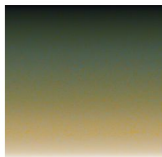
Figure 1. A self-evolving LLM plays three roles as PROPOSER, SOLVER and VERIFIER. The whole self-evolving process can be seen as different synthetic operations (synthesis qa, solution and feedback) on the same information source, which is the LLM itself.

Background on Epiplexity (Epistemic Complexity)

Random vs structural information

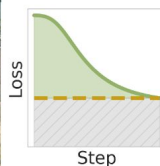
Low random info, low structural info

```
def is_even(n):
    if n == 0: return True
    elif n == 1: return False
    elif n == 2: return True
    elif n == 3: return False
    elif n == 4: return True
    elif n == 5: return False
    ...
```



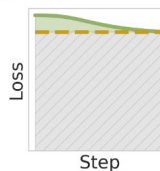
Moderate random info, high structural info

```
def dijkstra(g, s):
    D = defaultdict(lambda: float('inf'))
    D[s] = 0; q = [(0, s)]
    while q:
        d, u = pop(q)
        if d == D[u]:
            for v, w in g.get(u, []):
                if (nd := d + w) < D[v]:
                    D[v] = nd; push(q, (nd, v))
    return D
```



High random info, low structural info

```
API_KEY = "sk_7aF2jK1ycP9LmvYzz34"
USER_ID = "usr_4f8a2c1e9b7d3065"
BUCKET = "s3://data-8a3f1b-west-prod"
SAVE_DIR = "/mnt/marc/exp_7f2a/ckpts"
SAVE_CKPT = True
DEBUG = False
SEED = 9284715
...
```



Definition 8 (Epiplexity and Time-Bounded Entropy) Consider a random variable X on $\{0, 1\}^n$. Let

$$P^* = \arg \min_{P \in \mathcal{P}_T} \{ |P| + \mathbb{E}[\log 1/P(X)] \} \quad (3)$$

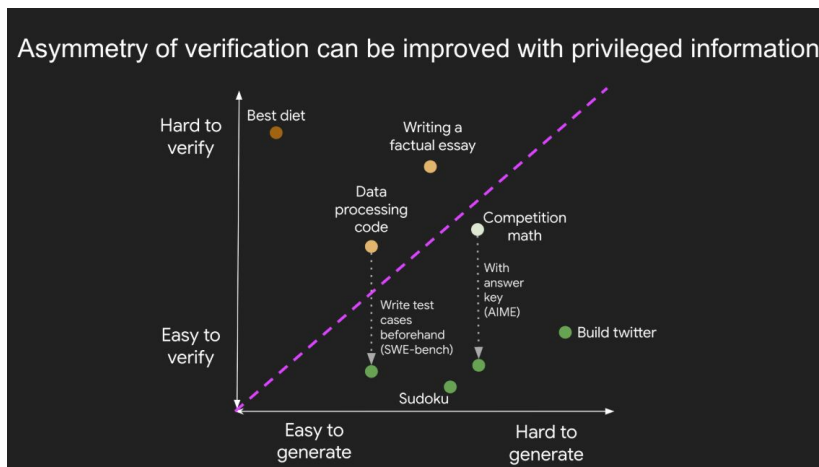
be the program that minimizes the time bounded MDL with ties broken by the smallest program, and expectations taken over X . $|P|$ denotes the length of the program P in bits, and logarithms are in base 2. We define the T -bounded epiplexity S_T and entropy H_T of the random variable X as

$$S_T(X) := |P^*|, \quad \text{and} \quad H_T(X) := \mathbb{E}[\log 1/P^*(X)]. \quad (4)$$

It explains something...
And help us find what is missing...

1. Why Self-Play Learns: Asymmetry

- Explain:
 - **Computation creates learnable information** → **different synthetic direction yield different amount of learnable information**
 - So it creates information asymmetric gap
 - This gap motivates the weak-to-strong self-training



1. Why Self-Play Learns: Asymmetry

- Explain:
 - **Computation creates learnable information** → **different synthetic direction yield different amount of learnable information**
 - So it creates information asymmetric gap
 - This gap motivates the weak-to-strong self-training

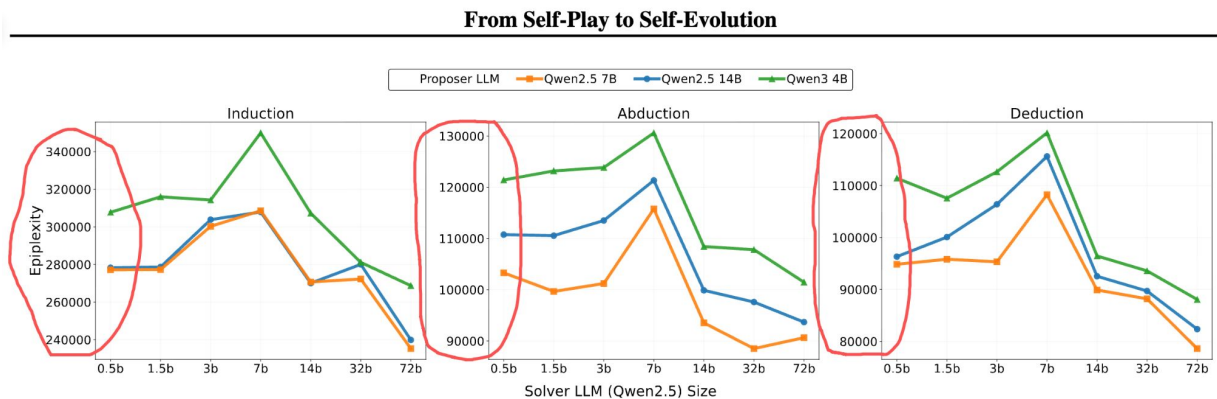


Figure 5. Epiplexity results on synthetic data with different tasks (induction, abduction and deduction) proposed by different PROPOSER LLMs and observed by different SOLVER LLMs. See details of calculating epiplexity in Appendix B.

1. Why Self-Play Learns: Asymmetry

- What is Missing?
 - Only weak verifier + proposer to strong solver, The ladder is not in closed loop
 - Only easy to verify → jagged intelligence
- Strong S → Strong V
 - Verifier-free RL
- Strong S → Strong P
 - Back-translation

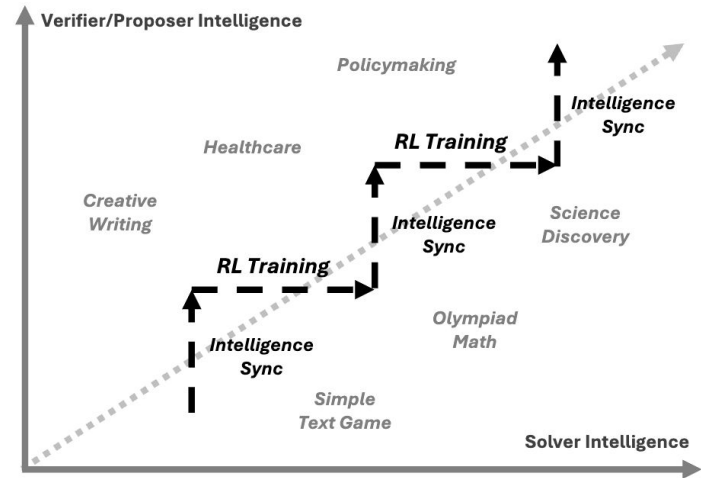


Figure 4. Climbing the intelligence asymmetry ladder by closing the loop among PROPOSER, SOLVER, and VERIFIER. “Intelligence synchronisation” denotes updating the weaker PROPOSER/VERIFIER with strong SOLVER. “Reinforcement learning” uses the weaker PROPOSER/VERIFIER to train the SOLVER.

Verifier-free RL: verifier co-evolve with solver

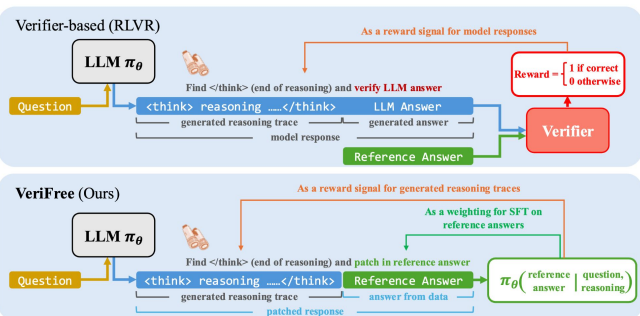


Figure 2: VeriFree enables RL-Zero-style LLM training without requiring access to a verifier. In the case of a single correct answer format, VeriFree optimizes exactly the same objective as RL-Zero with a lower variance gradient estimator.

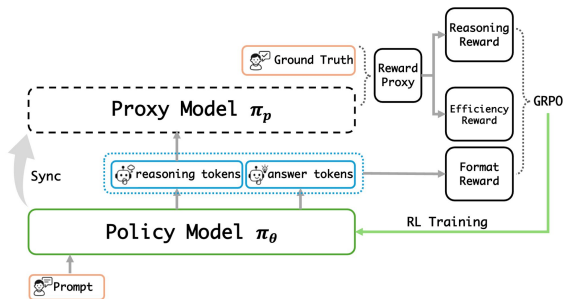
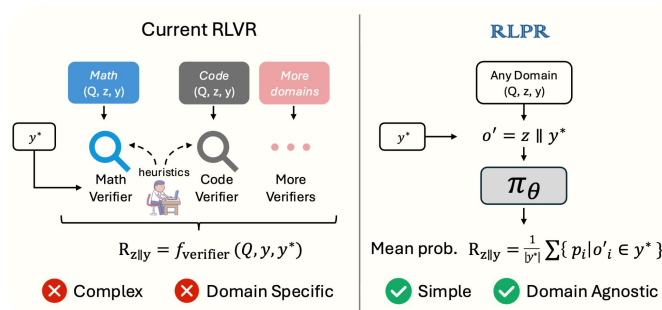
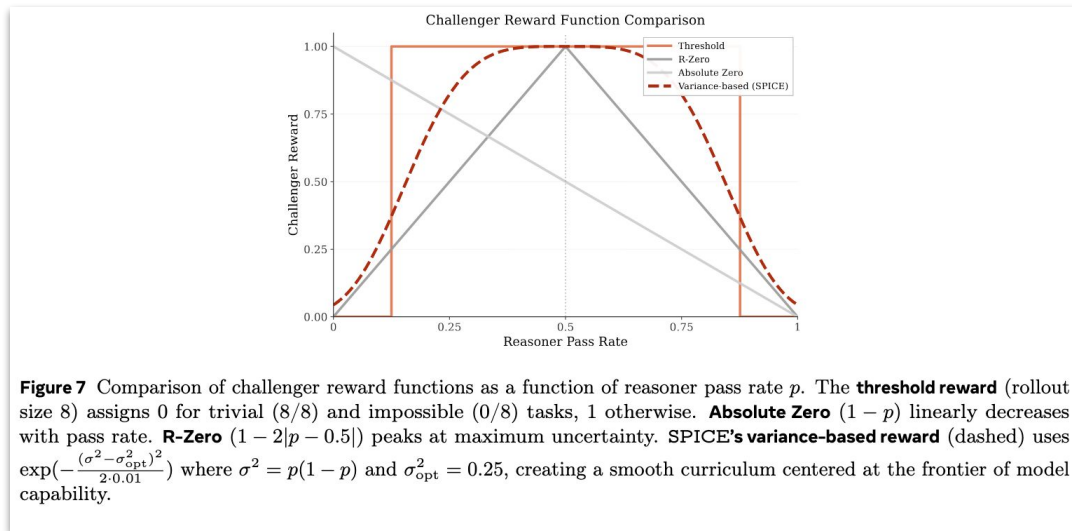


Figure 3: The overall process of NOVER.



- Three Verifier-free RL: VeriFree, NOVER, RLPR
- LLM acts as both the solver and verifier, **they share the same target**: maximize the probability to generate ground truth answer conditioned on its rollout reasoning parts

2. Why 50% acc for proposer?



- Proposer raise questions to challenge solver with 50% pass ratio

2. Why 50% acc for proposer: Capacity Match

- Explain in **capacity-data match**
 - Too easy for model, no emergent learnable structure
 - Too hard for model, just unlearnable noise

From Self-Play to Self-Evolution

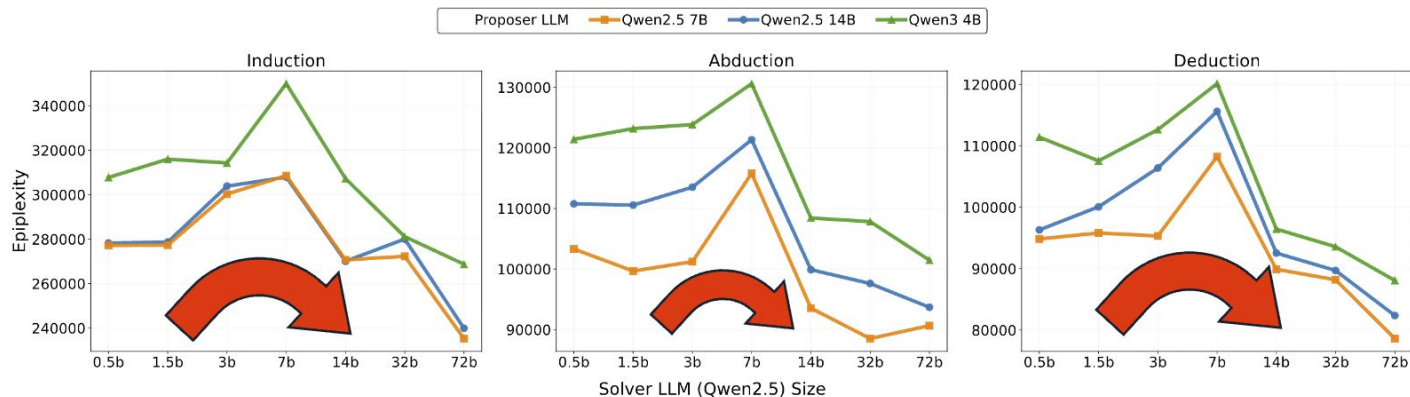


Figure 5. Epilexity results on synthetic data with different tasks (induction, abduction and deduction) proposed by different PROPOSER LLMs and observed by different SOLVER LLMs. See details of calculating epilexity in Appendix B.

2. What is Missing?

- Capacity Growth
 - **Information is $f(\text{data})$, Learnable information is $f(\text{data}, \text{model})$**
 - Capacity of model should grow with self-evolving iteration
 - Capacity means spent computation, factors including model size, test-time budget, activated/useful component, etc

3. Shannon said: of course we need external information

- But in what fashion?
 - Give **proactivity** to model: proposer knows what to challenge solver
 - Seek for **asymmetry**: not just more data
 - Env-Model **Co-Evolution**: the environment shapes how information is presented to the model

What Three Designs Actually Shapes

From Self-Play to Self-Evolution

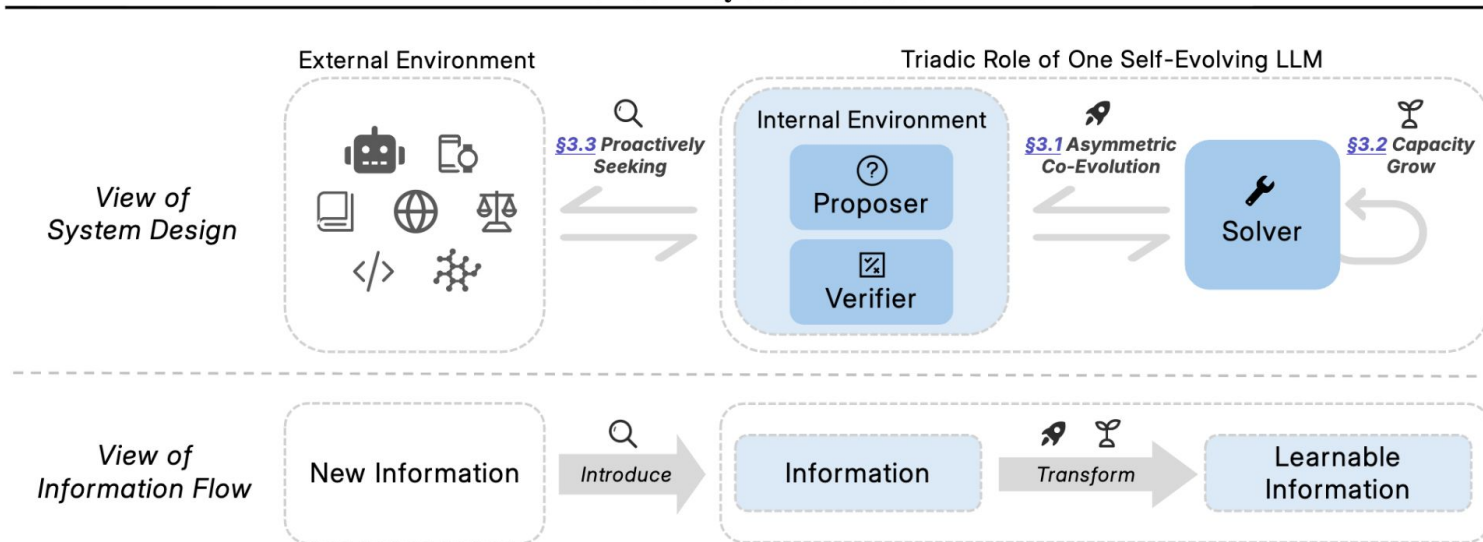


Figure 2. Overall framework of a triadic self-evolving loop. A self-evolving LLM plays three roles: the PROPOSER and VERIFIER form the internal environment, proactively interacting with the external environment to provide data and supervision for the SOLVER. The SOLVER and internal environment co-evolve asymmetrically, adaptively expanding capacity to capture more learnable information. From an information perspective, the system continually absorbs external information, and transform them into internal learnable information.

What Three Designs Actually Shapes

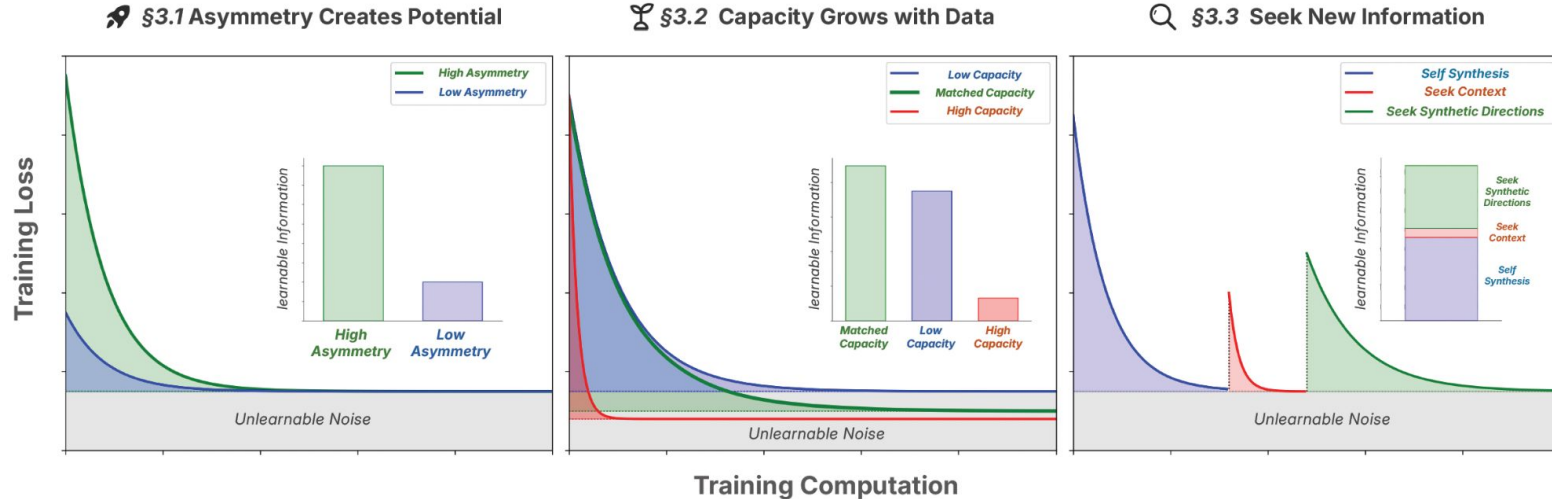
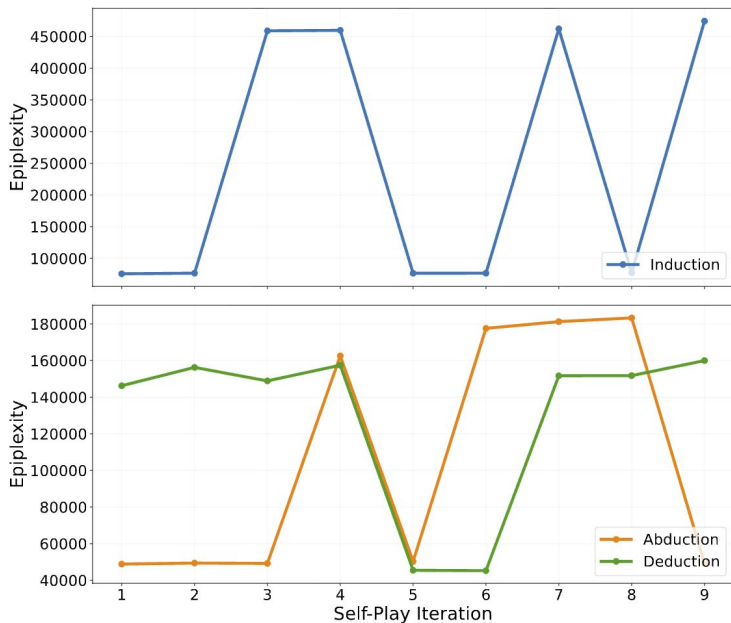


Figure 3. Illustration of three designs from the perspective of learnable information. Asymmetry between the SOLVER and the PROPOSER/VERIFIER creates learning opportunities. Expanding model capacity to match self evolving data opens space for learnable information. Reusing the same patterns in new contexts yields limited gains, whereas introducing new synthetic directions creates fresh asymmetries and thus new sources of learnable information.

Learnable Information Gain in Self-Play is Unstable



- Across the three coding tasks (**Induction**, **Abduction** and **Deduction**), the learnable information in self synthesized data varies substantially across self play iterations.

Figure 6. Epiplexity results during the self-play training on three tasks. See task detail in Appendix C.

Call to Action

- Sustainable Self-Evolving
 - three missing designs based on learnable information
- Beyond the Math&Code
 - cross the easy-to-verify side, Llm application in the wild
- Beyond the RL
 - is critic-free, model-free, policy-based RL all you need?
 - is countless-hyperparams learning algorithm all you need?
- Monitoring and Evaluation
 - from chaotic RL into auditable data engineering
 - clean and robust evaluation