

Archon-Prime: A Self-Reflective Meta-Reasoning Architecture for Neuro-Symbolic AI Systems

Anonymous Authors
Institution Name
email@institution.edu

July 9, 2025

Abstract

Contemporary AI systems exhibit a fundamental dichotomy: neural approaches offer flexibility but lack verifiability, while symbolic systems provide rigor but suffer from brittleness. We present **Archon-Prime**, a neuro-symbolic architecture that transcends this dichotomy through explicit meta-reasoning capabilities. Building upon the HAK-GAL (Hybrid AI Knowledge - Grounded Axiomatic Logic) framework, Archon-Prime implements a self-reflective reasoning system that not only operates over external domains but also maintains epistemological awareness of its own inferential processes. The architecture integrates three fundamental modes of reasoning—deduction, abduction, and induction—unified under Peirce’s theory of scientific inference. Its core innovation lies in a theory-informed, empirically adaptive portfolio manager that allocates computational resources based on a dual-layer complexity analysis combining theoretical complexity classes with machine-learned runtime predictions. Through an epistemic feedback loop utilizing a metamathematical ledger, the system achieves continuous self-optimization while maintaining formal guarantees within decidable fragments. We argue that Archon-Prime represents not a step toward artificial general intelligence, but rather an evolutionary synthesis of established techniques that creates a more robust, transparent, and scientifically grounded reasoning infrastructure—a system that knows what it knows and, crucially, knows what it does not know.

1 Introduction

The contemporary landscape of artificial intelligence is characterized by a fundamental tension between two paradigms. On one side, sub-symbolic approaches, exemplified by Large Language Models (LLMs), demonstrate remarkable flexibility and pattern recognition capabilities but operate as

epistemological black boxes, offering neither formal guarantees nor causal understanding [?]. On the other side, symbolic reasoning systems provide mathematical rigor and verifiability but struggle with scalability and exhibit brittleness when confronted with real-world complexity [?].

This dichotomy reflects a deeper philosophical divide that has persisted since the inception of AI: the tension between Connectionist and Symbolic approaches to intelligence [?]. Recent years have witnessed increasing interest in neuro-symbolic integration as a potential resolution [?], yet most approaches treat this integration at a superficial level, using neural networks merely as feature extractors for symbolic reasoners or employing symbolic constraints to guide neural training.

1.1 The Epistemic Gap in Modern AI

Current AI systems, regardless of their paradigm, share a critical limitation: they lack *epistemic self-awareness*. They cannot reliably assess their own knowledge boundaries, distinguish between degrees of certainty, or adapt their reasoning strategies based on problem complexity. This epistemic opacity manifests in several ways:

1. **The Black Box Problem:** Neural systems produce outputs without explanatory traces, making verification impossible.
2. **The Brittleness Problem:** Symbolic systems fail catastrophically outside their designed domains.
3. **The Scalability Problem:** Neither approach gracefully handles the complexity gradient from trivial to undecidable problems.
4. **The Learning Problem:** Systems either learn from data (neural) or from axioms (symbolic), but rarely both.

1.2 Our Contribution

We present **Archon-Prime**, a meta-reasoning architecture that addresses these limitations through explicit self-reflection and multi-modal reasoning integration. Our contributions are:

1. A **unified reasoning framework** integrating deduction, abduction, and induction based on Peirce’s theory of scientific inference.
2. A **dual-layer complexity analyzer** combining theoretical complexity classification with empirical runtime prediction.
3. An **adaptive prover portfolio manager** treating theorem proving as a multi-armed bandit problem.

4. An **epistemic self-improvement loop** utilizing a metamathematical ledger for continuous optimization.
5. A **formal analysis** of the system’s theoretical guarantees and fundamental limitations.

2 Theoretical Foundations

2.1 Peirce’s Triadic Theory of Inference

Charles Sanders Peirce identified three fundamental modes of logical inference that together constitute scientific reasoning [?]:

Definition 1 (Deduction). *Given premises P and rule $P \rightarrow Q$, infer conclusion Q . This is truth-preserving and forms the basis of mathematical proof.*

Definition 2 (Abduction). *Given observation Q and rule $P \rightarrow Q$, hypothesize P as a plausible explanation. This is hypothesis generation and forms the basis of scientific discovery.*

Definition 3 (Induction). *Given multiple instances of $P_i \rightarrow Q_i$, generalize to rule $\forall x. P(x) \rightarrow Q(x)$. This is pattern extraction and forms the basis of learning.*

Traditional reasoning systems implement only deduction, leaving critical gaps in their ability to handle incomplete knowledge or generate new hypotheses.

2.2 Computational Complexity as Foundational Constraint

The theoretical limits of computation impose fundamental constraints on any reasoning system:

Theorem 4 (Undecidability of FOL). *First-order logic is semi-decidable: there exists an algorithm that will find a proof if one exists, but no algorithm can definitively determine non-provability in finite time.*

This result, combined with Rice’s theorem on the undecidability of non-trivial semantic properties, establishes hard boundaries on what any reasoning system can achieve:

Theorem 5 (Rice’s Theorem). *Let \mathcal{P} be any non-trivial property of partial computable functions. Then $\{i : \phi_i \text{ has property } \mathcal{P}\}$ is undecidable.*

These theoretical limits necessitate a system design that explicitly acknowledges and manages undecidability rather than attempting to circumvent it.

2.3 The Polynomial Hierarchy and Fragment Classification

Different logical fragments exhibit different computational complexity:

- **Propositional Logic:** NP-complete for satisfiability
- **Quantified Boolean Formulas:** PSPACE-complete
- **First-Order Logic:** Semi-decidable (RE-complete)
- **Decidable Fragments:** Including description logics, guarded fragments

This hierarchy informs our approach to complexity analysis and prover selection.

2.4 Causal Reasoning and the Do-Calculus

Following Pearl’s causal hierarchy [?], we distinguish between:

1. **Association:** $P(Y|X)$ - observational relationships
2. **Intervention:** $P(Y|\text{do}(X))$ - causal effects
3. **Counterfactuals:** $P(Y_x|X', Y')$ - hypothetical reasoning

This framework enables reasoning about system interventions and hypothetical scenarios.

3 The Archon-Prime Architecture

3.1 System Overview

Archon-Prime consists of four interconnected layers, each addressing specific aspects of the meta-reasoning challenge:

3.2 The Multi-Modal Reasoning Core

The reasoning core unifies three inference modes within a single SMT-based framework:

3.2.1 Deductive Reasoning

Traditional theorem proving: Given knowledge base KB and query Q , determine if $KB \models Q$.

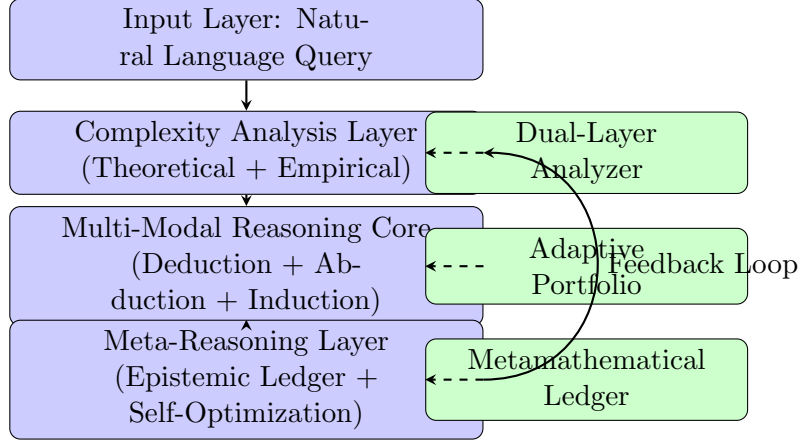


Figure 1: High-level architecture of Archon-Prime showing the four main layers and the epistemic feedback loop

Algorithm 1 Deductive Reasoning

Require: Knowledge base KB , Query Q

Ensure: Proof result (*proven, certificate*)

```

1:  $\phi_{KB} \leftarrow \text{encode\_to\_SMT}(KB)$ 
2:  $\phi_Q \leftarrow \text{encode\_to\_SMT}(Q)$ 
3:  $\text{solver} \leftarrow \text{create\_SMT\_solver}()$ 
4:  $\text{solver.add}(\phi_{KB})$ 
5:  $\text{solver.add}(\neg\phi_Q)$ 
6: if  $\text{solver.check}() = \text{UNSAT}$  then
7:   return (true, solver.get_proof())
8: else
9:   return (false, solver.get_model())
10: end if

```

3.2.2 Abductive Reasoning

Hypothesis generation: Given KB and observation Q where $KB \not\models Q$, find minimal H such that $KB \wedge H \models Q$.

Definition 6 (Minimal Abductive Explanation). H is a minimal abductive explanation for Q given KB if:

1. $KB \wedge H \models Q$ (*explanatory power*)
2. $KB \wedge H \not\models \perp$ (*consistency*)
3. $\nexists H' \subset H : KB \wedge H' \models Q$ (*minimality*)

Algorithm 2 Rigorous Abduction Engine

Require: Knowledge base KB , Goal Q , Maximum literals k

Ensure: Minimal hypothesis H or \emptyset

```
1: for  $size = 1$  to  $k$  do
2:    $H_{template} \leftarrow \text{create\_hypothesis\_variables}(size)$ 
3:    $solver \leftarrow \text{create\_SMT\_solver}()$ 
4:    $solver.add(KB \wedge H_{template} \models Q)$ 
5:    $solver.add(\text{consistency\_constraints}(H_{template}))$ 
6:   if  $solver.check() = \text{SAT}$  then
7:      $model \leftarrow solver.get\_model()$ 
8:     return  $\text{extract\_hypothesis}(model, H_{template})$ 
9:   end if
10: end for
11: return  $\emptyset$ 
```

3.2.3 Inductive Reasoning

Pattern generalization: Given observations $O = \{(x_i, y_i)\}$, find rule R such that $\forall i : R(x_i) = y_i$.

This component interfaces with Inductive Logic Programming (ILP) systems to learn horn clauses from examples.

3.3 The Dual-Layer Complexity Analyzer

Complexity analysis operates at two levels:

3.3.1 Theoretical Complexity Classification

Algorithm 3 Theoretical Complexity Analysis

Require: Formula ϕ in parsed AST form

Ensure: Complexity class \mathcal{C}

```
1:  $alternations \leftarrow \text{count\_quantifier\_alternations}(\phi)$ 
2:  $fragment \leftarrow \text{identify\_logical\_fragment}(\phi)$ 
3: if  $fragment = \text{PROPOSITIONAL}$  then
4:   return NP
5: else if  $fragment = \text{QBF} \wedge alternations \leq k$  then
6:   return  $\Sigma_k^P$  or  $\Pi_k^P$ 
7: else if  $fragment \in \text{DECIDABLE\_FRAGMENTS}$  then
8:   return  $\text{complexity\_of\_fragment}(fragment)$ 
9: else
10:  return RE (recursively enumerable)
11: end if
```

3.3.2 Empirical Runtime Prediction

A gradient-boosted tree model trained on historical proof attempts:

$$\hat{t} = f(\text{features}(\phi), \text{prover_type}, \text{domain_characteristics}) \quad (1)$$

Features include:

- Syntactic: clause count, literal count, variable count
- Structural: clause graph metrics (treewidth, connectivity)
- Semantic: predicate arities, function nesting depth
- Historical: success rate on similar formulas

3.4 The Adaptive Prover Portfolio Manager

Portfolio selection is formulated as a multi-armed bandit problem:

Definition 7 (Prover Selection as MAB). *Given:*

- Set of provers $\mathcal{P} = \{P_1, \dots, P_n\}$
- History $H = \{(\phi_i, P_j, \text{success}_i, \text{time}_i)\}$
- Current formula ϕ

Select subset $S \subseteq \mathcal{P}$ maximizing expected utility:

$$U(S, \phi) = \sum_{P \in S} p_{\text{success}}(P, \phi) \cdot v(\text{time}(P, \phi)) \quad (2)$$

where $v(\cdot)$ is a value function penalizing long runtimes.

The manager employs Thompson sampling for exploration-exploitation balance:

3.5 The Epistemic Self-Improvement Loop

The metamathematical ledger records all reasoning attempts:

Listing 1: Metamathematical Ledger Schema

```
@dataclass
class ReasoningRecord:
    timestamp: datetime
    formula: str
    complexity_profile: ComplexityReport
    portfolio_used: List[ProverInfo]
    result: ProofResult
```

Algorithm 4 Adaptive Portfolio Selection

Require: Formula ϕ , Complexity report C , Prover capabilities \mathcal{C}

Ensure: Portfolio S with resource allocations

```
1:  $eligible \leftarrow \{P \in \mathcal{P} : \text{can\_handle}(P, C.\text{fragment})\}$ 
2: for each  $P \in eligible$  do
3:    $\theta_P \sim \text{Beta}(\alpha_P, \beta_P)$  {Thompson sampling}
4:    $\hat{t}_P \leftarrow \text{predict\_runtime}(P, \phi)$ 
5:    $utility_P \leftarrow \theta_P / (1 + \hat{t}_P)$ 
6: end for
7:  $S \leftarrow \text{top\_k}(eligible, utility, k)$ 
8: for each  $P \in S$  do
9:    $budget_P \leftarrow \hat{t}_P \times \text{confidence\_factor}$ 
10: end for
11: return  $(S, \{budget_P\}_{P \in S})$ 
```

```
abductive_hypotheses: List[Hypothesis]
total_time: float
resource_usage: ResourceMetrics
```

This data drives three self-improvement mechanisms:

3.5.1 Model Retraining

Periodic retraining of the empirical runtime predictor:

$$\mathcal{L}_{new} = \mathcal{L}_{old} \cup \{(\phi_i, t_i^{actual})\}_{recent} \quad (3)$$

3.5.2 Portfolio Strategy Adaptation

Bayesian updates to prover success probabilities:

$$\alpha_P \leftarrow \alpha_P + \text{successes}_P, \quad \beta_P \leftarrow \beta_P + \text{failures}_P \quad (4)$$

3.5.3 Knowledge Base Evolution

High-confidence abductive hypotheses are candidates for KB extension:

$$KB_{new} = KB \cup \{H : \text{confidence}(H) > \tau \wedge \text{human_approved}(H)\} \quad (5)$$

4 Critical Analysis and Limitations

4.1 Theoretical Limitations

Despite its sophisticated architecture, Archon-Prime operates within fundamental theoretical constraints:

Theorem 8 (Incompleteness of Meta-Reasoning). *No formal system can completely capture its own reasoning processes without introducing paradox or incompleteness.*

Proof Sketch. By Tarski’s undefinability theorem, truth in a formal system cannot be defined within that system. Since meta-reasoning requires reasoning about the truth of reasoning processes, complete self-reflection is impossible within a single formal framework. \square

4.2 Practical Limitations

4.2.1 The Abduction Complexity Barrier

Proposition 9. *Finding minimal abductive explanations is NP-hard even for propositional logic.*

This necessitates heuristic approximations that may miss optimal hypotheses.

4.2.2 Distribution Shift in Self-Learning

The epistemic feedback loop assumes stationarity:

$$P(\text{formula}_{\text{future}}) \approx P(\text{formula}_{\text{past}}) \quad (6)$$

This assumption breaks down when problem domains evolve, potentially degrading performance.

4.2.3 The Human Bottleneck

Hypothesis validation requires human oversight, limiting the system’s autonomous improvement rate to human response time.

4.3 Comparison with Related Work

System	Deduction	Abduction	Meta-Reasoning	Self-Improvement
Traditional ATP	✓	×	×	×
Neural Theorem Provers	✓	×	×	Limited
Cognitive Architectures	✓	Limited	✓	×
Archon-Prime	✓	✓	✓	✓

Table 1: Comparison of reasoning capabilities across different system types

5 Implementation Considerations

5.1 Component Technologies

- **SMT Solver:** Z3 or CVC5 for core reasoning
- **ILP System:** Aleph or Metagol for induction
- **ML Framework:** XGBoost for runtime prediction
- **Database:** Graph database for metamathematical ledger
- **Interface:** REST API for human-in-the-loop validation

5.2 Scalability Strategies

1. **Proof Caching:** Memoization of proven subgoals
2. **Incremental Solving:** Reuse solver state across related queries
3. **Distributed Portfolio:** Parallel prover execution across nodes
4. **Ledger Pruning:** Periodic removal of low-value historical data

6 Future Directions

6.1 Theoretical Extensions

1. **Probabilistic Meta-Reasoning:** Extending to uncertain beliefs
2. **Higher-Order Self-Reflection:** Reasoning about reasoning about reasoning
3. **Temporal Meta-Logic:** Incorporating time into self-knowledge

6.2 Practical Enhancements

1. **Active Learning:** Strategic query generation for knowledge gaps
2. **Transfer Learning:** Adapting to new domains efficiently
3. **Explainable Abduction:** Natural language hypothesis explanations

7 Conclusion

Archon-Prime represents an evolutionary synthesis rather than a revolutionary breakthrough. By explicitly acknowledging and operationalizing the fundamental tensions in automated reasoning—between completeness and decidability, between theoretical guarantees and practical efficiency, between formal rigor and adaptive learning—it creates a more honest and ultimately more useful reasoning infrastructure.

The system’s value lies not in solving undecidable problems or achieving artificial general intelligence, but in creating a reasoning framework that:

- Knows its own limitations
- Adapts to its problem domain
- Integrates multiple reasoning modalities
- Continuously improves through self-reflection

This is the state of the art as it should be—not as it is. Archon-Prime provides a blueprint for building reasoning systems that embrace, rather than hide, their epistemic boundaries. In doing so, it points toward a future where AI systems are not black boxes of purported intelligence, but transparent partners in the pursuit of knowledge.

References

A Formal Specifications

A.1 Logic Fragment Hierarchy

$$\text{PROP} \subset \text{HORN} \subset \text{DATALOG} \subset \text{DL} \subset \text{FOL} \subset \text{HOL} \quad (7)$$

A.2 Complexity Class Relationships

$$\text{P} \subseteq \text{NP} \subseteq \Sigma_2^P \subseteq \dots \subseteq \text{PSPACE} \subseteq \text{RE} \quad (8)$$

B Pseudocode for Core Algorithms

B.1 Meta-Reasoning Cycle

Algorithm 5 Complete Meta-Reasoning Cycle

Require: Query Q , Knowledge base KB

Ensure: Result with meta-information

```
1:  $complexity \leftarrow \text{analyze\_complexity}(Q)$ 
2:  $portfolio \leftarrow \text{select\_portfolio}(Q, complexity)$ 
3:  $result \leftarrow \text{parallel\_prove}(portfolio, KB, Q)$ 
4: if  $result.success$  then
5:    $\text{record\_success}(Q, result)$ 
6: else
7:    $hypothesis \leftarrow \text{abduce\_explanation}(KB, Q)$ 
8:    $\text{record\_failure}(Q, hypothesis)$ 
9: end if
10:  $\text{update\_models}()$  {Periodic}
11: return  $(result, hypothesis, complexity)$ 
```
