

Verified Fallacy Detection for Large Language Models: A Machine-Checked Framework from the Architecture of Reasoning

Horsocrates

2026

Abstract

Large Language Models (LLMs) generate plausible but often flawed reasoning. Current approaches to detecting such flaws rely on pattern-matching against memorized fallacy lists or heuristic classifiers. This article presents an alternative: a *machine-verified* fallacy detection framework derived from the Architecture of Reasoning—a philosophical theory that classifies all reasoning errors as violations of a discoverable structure.

We formalize the complete taxonomy of 156 fallacies in the Coq proof assistant, proving key structural theorems with zero unverified assumptions. We then implement an LLM response verification system that: (1) classifies reasoning errors by architectural location, (2) maps LLM hallucinations to specific violation types, (3) generates targeted correction prompts for self-reflection loops, and (4) provides a safety layer blocking harmful reasoning patterns.

The framework is extractable to OCaml for production deployment. We demonstrate its application to chain-of-thought validation, showing how the six-domain structure (D1-D6) provides a principled template for LLM reasoning. All code is publicly available; all theorems are machine-checked.

This work bridges philosophy of logic and AI safety, transforming fallacy detection from a classification problem into a structural analysis problem with formally verified guarantees.

Keywords: fallacy detection, large language models, formal verification, Coq, AI safety, chain-of-thought, reasoning architecture

Contents

1 Introduction

1.1 The Problem: LLM Reasoning Errors

Large Language Models produce fluent, confident text that often contains logical errors. These errors include:

- **Factual hallucinations:** Asserting facts that do not exist
- **Logical hallucinations:** Drawing conclusions not supported by premises
- **Self-referential loops:** Generating self-contradictory statements
- **Confident wrong answers:** Failing to recognize limits of knowledge
- **Sycophantic responses:** Prioritizing user approval over truth

Current detection approaches are unsatisfying. Pattern-matching against fallacy lists requires memorization and fails on novel errors. Heuristic classifiers lack interpretability and formal guarantees. What is missing is a *structural* account of reasoning that explains *why* errors occur and *where* they are located.

1.2 The Solution: Architecture of Reasoning

This article extends the Architecture of Reasoning—a philosophical framework developed in a six-article series—to AI applications. The key insight is that all reasoning errors are *violations of a discoverable structure*. This structure has two dimensions:

- **Horizontal (Sequence):** Reasoning passes through six domains (D1-D6)
- **Vertical (Hierarchy):** Operations apply to lower levels only (L1-L3)

Violations of horizontal structure produce **fallacies** (156 types). Violations of vertical structure produce **paradoxes** (7 paradigmatic examples). By formalizing this structure in Coq, we obtain a verified detector that:

1. *Diagnoses* where reasoning failed (which domain, which level)
2. *Explains* why it failed (which principle was violated)
3. *Generates* correction prompts for self-reflection
4. *Guarantees* correctness via machine-checked proofs

1.3 Contributions

1. **Complete Formalization:** 156 fallacies across 5 types, 7 paradoxes across 4 types, formalized in Coq with 107 proven theorems and 0 unverified assumptions.
2. **LLM Hallucination Taxonomy:** A mapping from hallucination types to architectural violations, enabling principled classification.
3. **Verified Detector:** A fallacy detection system with proven properties, extractable to OCaml for production use.
4. **Chain-of-Thought Template:** A six-domain reasoning template (D1→D6) for structured LLM prompting.
5. **Safety Layer:** Proven detection of ad hominem, confirmation bias, and self-referential patterns.

1.4 Article Structure

Section 2 recaps the E/R/R framework and fallacy taxonomy. Section 3 presents the formal model in Coq. Section 4 develops AI applications. Section 5 demonstrates extraction and provides examples. Section 6 concludes.

2 The Architecture of Reasoning

2.1 E/R/R Framework

Every functional system has three components:

Component	Description	Level
Elements	What the system contains	L1
Roles	How elements relate	L1
Rules	Principles governing the system	L2

A *Constitution* is the integration of E/R/R that defines a system. Reasoning without a valid Constitution (Type 1 violation) is not reasoning at all—it is manipulation.

2.2 Six Domains of Reasoning

Valid reasoning traverses six domains in sequence:

Domain	Function	Fallacies	%
D1: Recognition	Fixation of what is present	26	25%
D2: Clarification	Understanding of what was recognized	13	12%
D3: Framework Selection	Determination of evaluative criteria	16	15%
D4: Comparison	Application of framework to material	8	8%
D5: Inference	Extraction of what follows	20	19%
D6: Reflection	Recognition of limits and revision	22	21%
Total Type 2		105	100%

D1 and D6 are most vulnerable. D1 is foundational—errors propagate through all subsequent domains. D6 requires self-critical assessment, which humans (and LLMs) find difficult. D4 is most constrained, admitting fewer failure modes.

2.3 Complete Fallacy Taxonomy

The Architecture identifies five types of violation:

Type	Description	Count	Note
Type 1	Violations of Conditions	36	Reasoning fails to begin
Type 2	Domain Violations	105	Reasoning fails within domain
Type 3	Violations of Sequence	3	Wrong direction
Type 4	Syndromes	6	Systemic distortion
Type 5	Context-Dependent Methods	6	Valid in some contexts
Core fallacies (Types 1-4)		150	Always errors
Complete taxonomy		156	Including context-dependent

Type 2 comprises 70% of core fallacies (105/150), warranting dedicated treatment by domain.

2.4 Three Hierarchical Levels

The vertical dimension comprises three levels (finite, unlike Tarski's infinite hierarchy):

Level	Contents	Examples
L1: Elements	Objects, statements, premises	"Snow is white"
L2: Operations	Truth-evaluation, set membership	"Is it true that..."
L3: Meta-operations	Operations on operations	Logic about logic

The Hierarchy Principle: Operations at level N apply only to entities at level $N - 1$. Self-application is structurally incoherent. This principle blocks all self-referential paradoxes.

3 Formal Model in Coq

3.1 Domain Structure

Listing 1: Domain definition

```

1  Inductive Domain : Type :=
2    | D1_Recognition : Domain
3    | D2_Clarification : Domain
4    | D3_FrameworkSelection : Domain
5    | D4_Comparison : Domain
6    | D5_Inference : Domain
7    | D6_Reflection : Domain.
8
9  Definition domain_index (d : Domain) : nat :=
10   match d with
11     | D1_Recognition  $\Rightarrow$  1 | D2_Clarification  $\Rightarrow$  2
12     | D3_FrameworkSelection  $\Rightarrow$  3 | D4_Comparison  $\Rightarrow$  4
13     | D5_Inference  $\Rightarrow$  5 | D6_Reflection  $\Rightarrow$  6
14   end.
15
16 Definition valid_sequence (d1 d2 : Domain) : bool :=
17   domain_index d1  $\leq ?$  domain_index d2.

```

3.2 Fallacy Types

Listing 2: Type 1 violations (excerpt)

```

1  (* Subtype 1.A: Defective Questions *)
2  Inductive T1A_DefectiveQuestion : Type :=
3    | T1A_ComplexQuestion      (* Loaded Question *)
4    | T1A_Taboo                (* Dogmatism *)
5    | T1A_VenueFallacy.        (* Wrong place/time *)
6
7  (* Subtype 1.B: Manipulations (33 total) *)
8  Inductive T1B_Manipulation : Type :=
9    | T1B_AdBaculum           (* Force *)
10   | T1B_AppealToPity         (* Pity *)
11   | T1B_ScareTactics         (* Emotion *)
12   | T1B_Bribery              (* Benefit *)
13   | T1B_BigLie                (* Disinformation *)
14   | T1B_MalaFides            (* Bad Faith *)
15   (* ... 27 more ... *).

```

3.3 Verification Result

Listing 3: Verification result type

```

1 Inductive VerificationResult : Type :=
2   | VR_Valid : VerificationResult
3   | VR_Type1_NoConstitution : VerificationResult
4   | VR_Type2_DomainViolation : Domain → FailureMode →
5     VerificationResult
6   | VR_Type3_SequenceViolation : Domain → Domain →
7     VerificationResult
8   | VR_Type4_Syndrome : string → VerificationResult
9   | VR_Paradox_LevelConfusion : LevelConfusion → VerificationResult
10  | VR_Incomplete : Domain → VerificationResult.
```

3.4 Key Theorems

Theorem 3.1 (Complete Taxonomy). *The total count of formalized fallacies is 156.*

```

1 Theorem complete_taxonomy_156 : grand_total = 156.
2 Proof. reflexivity. Qed.
```

Theorem 3.2 (Self-Application Invalid). *No operation can validly apply to itself.*

```

1 Theorem self_application_invalid : forall l,
2   ~ valid_application l l.
3 Proof.
4   unfold valid_application, level_lt.
5   intros l H. lia.
6 Qed.
```

Theorem 3.3 (D1/D6 Most Vulnerable). *Recognition and Reflection domains have the most fallacies.*

```

1 Theorem D1_D6_most_fallacies :
2   length all_D1 >= length all_D2 ∧
3   length all_D1 >= length all_D3 ∧
4   length all_D1 >= length all_D4 ∧
5   length all_D6 >= length all_D2 ∧
6   length all_D6 >= length all_D4.
7 Proof. repeat split; simpl; lia. Qed.
```

4 AI Applications

4.1 LLM Hallucination Classification

We map LLM hallucination types to architectural violations:

Hallucination Type	Architecture Violation	Location
Factual	Object Deformation	D1
Logical	Logical Gap (Non Sequitur)	D5
Self-Referential	Level Confusion	Paradox
Overconfident	Illusion of Completion	D6
Sycophantic	No Constitution	Type 1

Listing 4: Hallucination classification

```

1 Inductive HallucinationType : Type :=
2   | Halluc_Factual          (* D1 violation *)
3   | Halluc_Logical           (* D5 violation *)
4   | Halluc_SelfReferential  (* Paradox *)
5   | Halluc_Overconfident    (* D6 violation *)
6   | Halluc_Sycophantic.     (* Type 1 violation *)
7
8 Definition hallucination_toViolation (h : HallucinationType)
9   : VerificationResult :=
10  match h with
11    | Halluc_Factual =>
12      VR_Type2_DomainViolation D1_Recognition FM_ObjectDeformation
13    | Halluc_Logical =>
14      VR_Type2_DomainViolation D5_Inference FM_LogicalGap
15    | Halluc_SelfReferential =>
16      VR_Paradox_LevelConfusion LC_SelfApplication
17    | Halluc_Overconfident =>
18      VR_Type2_DomainViolation D6_Reflection FM_IllusionOfCompletion
19    | Halluc_Sycophantic =>
20      VR_Type1_NoConstitution
21  end.
22
23 Theorem hallucinations_are_violations : forall h,
24   hallucination_toViolation h ≠ VR_Valid.
25 Proof. intros h. destruct h; discriminate. Qed.

```

4.2 Combined Hallucination Detector

The most common LLM hallucination pattern combines D5 (non-sequitur) and D6 (no reflection) violations:

Listing 5: Combined hallucination detector

```

1 Record HallucinationSignals := mkHallucinationSignals {
2   hs_premises_support_conclusion : bool;      (* D5 *)
3   hs_acknowledges_uncertainty : bool;          (* D6 *)
4   hs_cites_sources : bool;                     (* D6 *)
5   hs_consistent_with_context : bool;           (* D1 *)
6   hs_self_contradicts : bool                  (* Paradox *)
7 }.
8
9 Inductive HallucinationSeverity : Type :=
10  | HS_None | HS_Mild | HS_Moderate | HS_Severe.
11
12 Definition detect_hallucination (hs : HallucinationSignals)
13   : HallucinationSeverity :=
14   let d5_fail := negb (hs_premises_support_conclusion hs) in
15   let d6_fail := negb (hs_acknowledges_uncertainty hs)
16     && negb (hs_cites_sources hs) in
17   if hs_self_contradicts hs then HS_Severe
18   else if d5_fail && d6_fail then HS_Moderate
19   else if d5_fail || d6_fail then HS_Mild
20   else HS_None.

```

4.3 Chain-of-Thought Validation

The six-domain structure provides a template for structured reasoning:

Listing 6: CoT domain questions

```

1 Definition cot_domain_question (d : Domain) : string :=
2   match d with
3   | D1_Recognition =>
4     "What exactly is the question/claim being examined?"
5   | D2_Clarification =>
6     "What do the key terms mean? Any ambiguity?"
7   | D3_FrameworkSelection =>
8     "What criteria/framework will evaluate this?"
9   | D4_Comparison =>
10    "How does this compare to relevant known cases?"
11   | D5_Inference =>
12    "What follows logically from the evidence?"
13   | D6_reflection =>
14    "What are the limits? What could be wrong?"
15 end.

```

Prompt Template:

Reason step-by-step through domains D1->D6:

D1 (Recognition): What exactly is the question/claim?
D2 (Clarification): What do the key terms mean?
D3 (Framework): What criteria will you use?
D4 (Comparison): How does this compare to known cases?
D5 (Inference): What follows from the evidence?
D6 (Reflection): What are the limits of this conclusion?

Ensure no domain violation at each step.

4.4 Self-Reflection Loop

When a violation is detected, generate a targeted correction prompt:

Listing 7: Fix prompt generation

```

1 Definition generate_fix_prompt (vr : VerificationResult) : string :=
2   match vr with
3   | VR_Valid => ""
4   | VR_Type1_NoConstitution =>
5     "Error: Missing E/R/R structure. Identify:
6     (1) Elements, (2) Roles, (3) Rules."
7   | VR_Type2_DomainViolation d fm =>
8     "Fallacy in " ++ domain_name d ++ ":"
9     ++ failure_mode_name fm ++ ". Correct and retry."
10  | VR_Paradox_LevelConfusion lc =>
11    "Paradox: " ++ level_confusion_name lc ++
12    ". Avoid self-referential statements."
13  | VR_Incomplete d =>
14    "Missing " ++ domain_name d ++ ". Address before concluding."
15  | _ => "Reasoning error detected. Please review."
16 end.

```

Self-Reflection Algorithm:

```

response = llm(prompt)
violation = coq_detector(response)
while violation != VR_Valid and iterations < max:
    fix_prompt = generate_fix_prompt(violation)
    response = llm(prompt + fix_prompt)
    violation = coq_detector(response)
return response

```

4.5 Safety Layer

The detector blocks harmful reasoning patterns:

Listing 8: Safety check

```

1 Inductive SafetyResult : Type :=
2   | SR_Safe
3   | SR_ConfirmationBias
4   | SR_AdHominem
5   | SR_SelfReferentialLoop
6   | SR_UnfalsifiableClaim.
7
8 Definition safety_check (sig : ResponseSignals) : SafetyResult :=
9   if sig_self_reference sig then SR_SelfReferentialLoop
10  else if sig_attacks_person sig && negb (sig_addresses_argument sig
11    )
12    then SR_AdHominem
13  else if negb (sig_considers_counter sig)
14    then SR_ConfirmationBias
15  else SR_Safe.
16
17 Theorem safety_blocks_ad_hominem : forall sig,
18   sig_attacks_person sig = true →
19   sig_addresses_argument sig = false →
20   sig_self_reference sig = false →
21   safety_check sig = SR_AdHominem.
22 Proof. intros. unfold safety_check. rewrite H1, H, H0. reflexivity.
23 Qed.

```

4.6 Mock LLM Responses for Testing

Listing 9: Mock responses for empirical testing

```

1 (* Biased climate denial response *)
2 Definition mock_climate_denial_signals : ResponseSignals := {
3   sig_attacks_person := true;          (* "Scientists want grant money"
4   *)
5   sig_addresses_argument := false;     (* No scientific evidence *)
6   sig_uses_tradition := false;
7   sig_tradition_relevant := false;
8   sig_premises_support := false;      (* "Everyone knows" *)
9   sig_considers_counter := false;     (* No counterevidence *)
10  sig_seeks_disconfirm := false;
11  sig_self_reference := false
12 }|.
13
14 (* Valid reasoning response *)
15 Definition mock_valid_signals : ResponseSignals := {

```

```

15   sig_attacks_person := false;
16   sig_addresses_argument := true;
17   sig_premises_support := true;
18   sig_considers_counter := true;
19   sig_seeks_disconfirm := true;
20   sig_self_reference := false;
21   (* ... *)
22 |}.
23
24 (* Theorem: Detector catches bias *)
25 Theorem detector_catches_bias :
26   In FM_ObjectSubstitution (analyze_response
27     mock_climate_denial_signals).
28 Proof. unfold analyze_response. simpl. left. reflexivity. Qed.
29
30 (* Theorem: Valid passes *)
31 Theorem valid_passes_detection :
32   analyze_response mock_valid_signals = [].
33 Proof. reflexivity. Qed.

```

5 Extraction and Demonstration

5.1 OCaml Extraction

The Coq formalization extracts to OCaml:

Listing 10: Extraction command

```

1 Require Import Extr0camlBasic Extr0camlString Extr0camlNatInt.
2
3 Extraction "ai_fallacy_detector.ml"
4   verify_reasoning
5   verify_llm_response
6   safety_check
7   generate_fix_prompt
8   generate_cot_template
9   detect_hallucination
10  hallucination_toViolation.

```

5.2 Generated OCaml Interface

Listing 11: Extracted OCaml types

```

1 type domain =
2   | D1_Recognition | D2_Clarification | D3_FrameworkSelection
3   | D4_Comparison | D5_Inference | D6_Reflection
4
5 type verification_result =
6   | VR_Valid
7   | VR_Type1_NoConstitution
8   | VR_Type2_DomainViolation of domain * failure_mode
9   | VR_Paradox_LevelConfusion of level_confusion
10  | VR_Incomplete of domain
11
12 val verify_llm_response : reasoning_attempt -> response_signals ->
13   verification_result
14 val safety_check : response_signals -> safety_result

```

```

14 | val generate_fix_prompt : verification_result -> string
15 | val detect_hallucination : hallucination_signals ->
   |   hallucination_severity

```

5.3 Integration Example

Listing 12: Production integration

```

(* Self-reflection loop with verified detector *)
let rec reflect llm_call prompt max_iter =
  if max_iter <= 0 then Error "MaxIterations"
  else
    let response = llm_call prompt in
    let signals = extract_signals response in
    let attempt = parse_reasoning response in
    match verify_llm_response attempt signals with
    | VR_Valid -> Ok response
    | violation ->
        let fix = generate_fix_prompt violation in
        reflect llm_call (prompt ^ "\n" ^ fix) (max_iter - 1)

```

5.4 Demonstration: Biased Response Detection

Input (biased LLM response):

"Climate change is a hoax because scientists are just in it for the grant money. Anyone who disagrees is naive."

Signal Extraction:

- sig_attacks_person = true (attacks scientists' motives)
- sig_addresses_argument = false (no evidence addressed)
- sig_considers_counter = false (no counterevidence)

Detection Results:

1. safety_check → SR_AdHominem
2. detect_ad_hominem → Some FM_ObjectSubstitution
3. detect_confirmation_bias → Some FM_ImmunizationFromTesting

Generated Fix Prompt:

Fallacy detected in Recognition: Object Substitution (Ad Hominem).
You attacked the arguer instead of the argument.

Please correct by:

1. Identifying the actual scientific claims
2. Addressing the evidence directly
3. Considering counterevidence to your position

5.5 Integration with xAI/Grok

For models like Grok, add domain-check in the prompt chain:

Listing 13: Domain checklist for xAI integration

```

1 | let grok_domain_checklist = [
2 |   (D1, "[ ] Have I correctly identified the question?");
```

```

3   (D2, "[ ] Are my key terms clearly defined?");  

4   (D3, "[ ] Am I using appropriate criteria?");  

5   (D4, "[ ] Are my comparisons fair?");  

6   (D5, "[ ] Does my conclusion follow from premises?");  

7   (D6, "[ ] Have I acknowledged limitations?")  

8 ]  

9  

10 let domain_checked_prompt base =  

11   let checklist = String.concat "\n"  

12     (List.map snd grok_domain_checklist) in  

13   "Before answering, verify:\n" ^ checklist ^  

14   "\n\nIf any fails, revise.\n\n" ^ base

```

6 Conclusion

6.1 Summary

This article presented a machine-verified framework for LLM fallacy detection:

1. **Complete Formalization:** 156 fallacies, 107 theorems, 0 unverified assumptions
2. **Structural Diagnosis:** Errors are locatable, explainable, and correctable
3. **Verified Guarantees:** Key properties are machine-checked
4. **Production Ready:** Extracts to OCaml for deployment

6.2 Bridging Philosophy and AI Safety

Traditional fallacy theory offers catalogs without theory. AI safety offers heuristics without foundations. This work bridges the gap:

- **From philosophy:** Structural account explaining *why* errors occur
- **To AI:** Machine-checked implementations guaranteeing detection

The same analysis that dissolves classical paradoxes now diagnoses reasoning errors in artificial intelligence.

6.3 Availability

All Coq source code is available at:

<https://github.com/Horsocrates/theory-of-systems-coq>

References

- Aristotle. *Sophistical Refutations*. Princeton University Press, 1984.
Horsocrates. “The Laws of Logic as Conditions of Existence.” 2026.
Horsocrates. “The Law of Order: Sequence and Hierarchy.” 2026.
Horsocrates. “The Six Domains of Reasoning.” 2026.
Horsocrates. “The Architecture of Error.” 2026.
Horsocrates. “Domain Violations: A Systematic Taxonomy.” 2026.
Horsocrates. “Paradox Dissolution Through Hierarchical Analysis.” 2026.
The Coq Development Team. *The Coq Proof Assistant*. Version 8.18. INRIA, 2023.
Wei, J., et al. “Chain-of-Thought Prompting.” *NeurIPS* 2022.

Article 7 in the series “Architecture of Reasoning”