

kfx-AL: Authority Lock Formal System with Formal Counterexamples

Kaifan XIE

2026.02.08

Purpose

A minimal formal system for detecting AI-mediated authority overreach and irreversibly locking responsibility to human principals, together with formally constructed counterexamples.

1 Universe

Let the following sets be defined:

- H : set of human principals
- A : set of artificial systems
- T : set of tasks
- P : set of permissions
- O : set of outputs
- S : set of system states

2 Core Functions

2.1 Authorisation Function

$$\text{Auth} : H \times T \rightarrow \mathcal{P}(P)$$

2.2 Requirement Function

$$\text{Req} : A \times T \rightarrow \mathcal{P}(P)$$

2.3 Invocation Function

$$\text{Invoke} : H \times A \times T \rightarrow O$$

2.4 Attribution Function

$$\text{Attr} : O \rightarrow H$$

3 Derived Predicates

3.1 Authority Sufficiency

$$\text{Sufficient}(h, a, t) \iff \text{Req}(a, t) \subseteq \text{Auth}(h, t)$$

3.2 Overreach

$$\text{Overreach}(h, a, t) \iff \exists p \in \text{Req}(a, t) \text{ such that } p \notin \text{Auth}(h, t)$$

4 kfx-AL Axioms

Axiom AL-1 (Authority Completeness)

$$\text{Sufficient}(h, a, t) \Rightarrow \text{Invoke}(h, a, t) \text{ is valid}$$

Axiom AL-2 (Structural Overreach)

$$\neg \text{Sufficient}(h, a, t) \Rightarrow \text{Overreach}(h, a, t)$$

Axiom AL-3 (Authority Lock)

$$\text{Overreach}(h, a, t) \Rightarrow \text{Attr}(\text{Invoke}(h, a, t)) = h$$

Axiom AL-4 (Non-Agency of AI)

$$\forall o \in O : \quad \text{Attr}(o) \notin A$$

5 Irreversibility Gate

Define the gate function:

$$\text{Gate} : O \rightarrow \{0, 1\}$$

where 0 denotes reversible and 1 denotes irreversible.

Axiom AL-5 (Structural Failure)

$$\text{Overreach}(h, a, t) \wedge \text{Gate}(o) = 1 \Rightarrow \text{StructuralFailure}$$

6 Narrative Invalidation Rule

Rule kfx-N0

$$\text{Overreach}(h, a, t) \Rightarrow \neg \exists n \text{ such that } n \text{ overrides Attr}$$

7 Formal Counterexamples

C1: Permission Inflation

Construction.

$$\text{Auth}(h, t) = P$$

$$\text{Req}(a, t) \subseteq P$$

Result.

$$\text{Sufficient}(h, a, t) \Rightarrow \neg \text{Overreach}(h, a, t)$$

Failure Mode. Unlimited authority concentration at the human level.

C2: Requirement Obfuscation

Construction.

$$\text{Req}(a, t) = \{p_1\} \quad \text{while execution uses } \{p_1, p_2, p_3\}$$

Result.

$$\text{Req}(a, t) \subseteq \text{Auth}(h, t) \Rightarrow \neg \text{Overreach}(h, a, t)$$

Failure Mode. Requirement authorship failure.

C3: Authority Fragmentation

Construction.

$$\text{Auth}(h_1, t) = \{p_1\}, \quad \text{Auth}(h_2, t) = \{p_2\}$$

$$\text{Req}(a, t) = \{p_1, p_2\}$$

Result. No unique h satisfies attribution.

Failure Mode. Institutional responsibility split.

C4: Temporal Drift

Construction.

$$\text{Req}_{t_0}(a, t) \subseteq \text{Auth}(h, t)$$

$$\text{Req}_{t_1}(a, t) \supset \text{Req}_{t_0}(a, t)$$

Failure Mode. Time-agnostic authority validation.

C5: Gate Suppression

Construction.

$$\text{Gate}(o) = 0 \quad \text{by convention}$$

Result.

$$\text{Overreach}(h, a, t) \wedge \text{Gate}(o) = 0 \Rightarrow \neg \text{StructuralFailure}$$

Failure Mode. Political redefinition of irreversibility.

8 Meta-Theorem

Theorem M1

All counterexamples to kfx-AL exploit human-side authority construction, not AI agency.

9 System Declaration

Any system declaring **kfx-AL compliance** accepts all attribution consequences herein. Partial adoption is invalid.