

# Semantic Compression via Type Systems: Matroid Structure and Kolmogorov-Optimal Witnesses

Anonymous

January 15, 2026

## Abstract

We study semantic inference under observational constraints. An interface-only observer queries membership in a fixed family of interfaces  $\Phi_{\mathcal{I}} = \{q_I : I \in \mathcal{I}\}$  and must decide semantic properties such as type identity. We prove an information barrier: every interface-only observer is constant on indistinguishability classes (values with identical interface profiles), hence cannot compute any property that varies within such a class. In contrast, a nominal-tag observer—with access to a single type identifier per value—achieves constant witness cost:  $W(\text{type-identity}) = O(1)$  primitive queries. We further establish that minimal complete axis sets form the bases of a matroid, and that nominal-tag observers achieve the unique Pareto-optimal point in the  $(L, W, D)$  tradeoff (tag length, witness cost, distortion). All results are machine-checked in Lean 4.

**Keywords:** semantic inference, observational constraints, information barriers, witness cost, matroid structure, type systems, Lean 4

## 1 Introduction

### 1.1 Observational Constraints and Semantic Inference

Consider the following inference problem: a procedure observes a program value and must determine its semantic properties (e.g., type identity, provenance). The procedure’s access is restricted to a fixed family of *interface observations*—predicates that test membership in declared interfaces.

**Definition 1.1** (Interface observation family). Fix a set of interfaces  $\mathcal{I}$ . For each  $I \in \mathcal{I}$ , define the interface-membership observation  $q_I : \mathcal{V} \rightarrow \{0, 1\}$ , where  $q_I(v) = 1$  iff  $v$  satisfies interface  $I$ . Let  $\Phi_{\mathcal{I}} = \{q_I : I \in \mathcal{I}\}$ .

**Definition 1.2** (Interface profile). Define  $\pi : \mathcal{V} \rightarrow \{0, 1\}^{\mathcal{I}}$  by  $\pi(v) = (q_I(v))_{I \in \mathcal{I}}$ .

**Definition 1.3** (Interface indistinguishability).  $v \sim w$  iff  $\pi(v) = \pi(w)$ .

**Definition 1.4** (Interface-only procedure). An interface-only procedure is any algorithm whose interaction with a value is limited to queries in  $\Phi_{\mathcal{I}}$ .

The central question is: **what semantic properties can an interface-only procedure compute?**

## 1.2 The Impossibility Barrier

**Theorem 1.5** (Information barrier from interface-only evidence). *Every interface-only procedure is constant on  $\sim$ -equivalence classes. Consequently, no interface-only procedure can compute any property that differs for some  $v \sim w$ .*

This is an information barrier: the restriction is not computational (unbounded time/memory does not help) but informational (the evidence itself is insufficient).

## 1.3 The Positive Result: Nominal Tagging

In contrast, nominal tagging—storing an explicit type identifier per value—provides constant-size evidence for type identity.

**Definition 1.6** (Witness cost). Let  $W(P)$  denote the minimum number of primitive queries required to compute property  $P$ . A primitive query is either an interface observation  $q_I \in \Phi_{\mathcal{I}}$  or a nominal-tag access (reading the type identifier).

**Theorem 1.7** (Constant witness for nominal type identity). *Nominal-tag access admits a constant-cost witness for type identity:  $W(\text{type-identity}) = O(1)$  primitive queries.*

## 1.4 Main Contributions

1. **Impossibility Theorem:** No interface-only procedure can compute properties that vary within  $\sim$ -equivalence classes (Theorem 1).
2. **Constant-Witness Result:** Nominal tagging achieves  $W(\text{type-identity}) = O(1)$  (Theorem 2).
3. **Equicardinality Theorem:** All minimal complete type axis sets have equal cardinality (a consequence of matroid-like structure).
4. **Rate–Witness–Distortion Optimality:** Nominal-tag observers achieve the unique Pareto-optimal point in the  $(L, W, D)$  tradeoff (tag length, witness cost, distortion).
5. **Machine-Checked Proofs:** All results formalized in Lean 4 ( $\sim 6,000$  lines, 265 theorems, 0 sorry).

## 1.5 Audience and Scope

This paper is written for the information theory and compression community. We assume familiarity with matroid theory and basic information-theoretic concepts. We provide concrete instantiations in widely used programming language runtimes (CPython, Java, TypeScript, Rust) as corollaries to the main theorems.

# 2 Compression Framework

## 2.1 Formal Model: Observations and Equivalence

Let  $\mathcal{V}$  denote the space of all program values,  $\mathcal{I}$  the set of interfaces, and  $\Phi_{\mathcal{I}}$  the interface observation family (Definition 1).

**Definition 2.1** (Interface equivalence). Values  $v, w \in \mathcal{V}$  are interface-equivalent, written  $v \sim w$ , iff  $\pi(v) = \pi(w)$ —i.e., they satisfy exactly the same interfaces.

An interface-only procedure can only distinguish values that are not interface-equivalent. Therefore, any property computed by an interface-only procedure must be constant on  $\sim$ -equivalence classes.

## 2.2 Witness Cost

A *witness* for a property  $P$  is a procedure that, given access to a value, computes  $P$  using primitive queries.

**Definition 2.2** (Witness cost). The witness cost of property  $P$  is  $W(P) = \min\{c(w) : w \text{ is a witness procedure for } P\}$ , where  $c(w)$  is the number of primitive queries (interface observations or nominal-tag accesses) required by  $w$ .

*Remark 2.3* (Connection to algorithmic information theory). Witness cost is related to Kolmogorov complexity, but measures query count under a fixed primitive set rather than description length under a universal machine. This makes  $W$  a concrete, computable quantity suitable for comparing practical systems.

## 2.3 Rate–Witness–Distortion Tradeoff

We analyze type identity checking under three dimensions:

**Definition 2.4** (Tag length). The tag length  $L$  is the number of machine words required to store a type identifier per value. (Under a fixed word size  $w$ , this corresponds to  $\Theta(w)$  bits.)

**Definition 2.5** (Witness cost). The witness cost  $W$  is the minimum number of primitive queries required to implement type identity checking (Definition above).

**Definition 2.6** (Distortion). The distortion  $D$  is a worst-case semantic failure flag:

$$D = 0 \iff \forall v_1, v_2 [\text{type}(v_1) = \text{type}(v_2) \Rightarrow \text{behavior}(v_1) \equiv \text{behavior}(v_2)]$$

Otherwise  $D = 1$ . Here  $\text{behavior}(v)$  denotes the observable behavior of  $v$  under program execution (e.g., method dispatch outcomes).

A type system is characterized by a point  $(L, W, D)$  in this three-dimensional space. The question is: which points are achievable, and which are Pareto-optimal?

## 3 Matroid Structure

### 3.1 Type Axes

A *type axis* is a semantic dimension along which types can vary. Examples:

- **Identity:** Explicit type name or object ID
- **Structure:** Field names and types
- **Behavior:** Available methods and their signatures

- **Scope:** Where the type is defined (module, package)
- **Mutability:** Whether instances can be modified

A *complete* axis set distinguishes all semantically distinct types. A *minimal complete* axis set is complete with no proper complete subset.

### 3.2 Matroid Structure of Type Axes

**Definition 3.1** (Axis bases family). Let  $E$  be the set of all type axes. Let  $\mathcal{B} \subseteq 2^E$  be the family of minimal complete axis sets.

**Lemma 3.2** (Basis exchange). *For any  $B_1, B_2 \in \mathcal{B}$  and any  $e \in B_1 \setminus B_2$ , there exists  $f \in B_2 \setminus B_1$  such that  $(B_1 \setminus \{e\}) \cup \{f\} \in \mathcal{B}$ .*

*Proof.* See Lean formalization: `proofs/axis_framework.lean`, lemma `basis_exchange`. ■

**Theorem 3.3** (Matroid bases).  *$\mathcal{B}$  is the set of bases of a matroid on ground set  $E$ .*

*Proof.* By the basis-exchange lemma and the standard characterization of matroid bases. ■

**Corollary 3.4** (Well-defined semantic dimension). *All minimal complete axis sets have equal cardinality. Hence the “semantic dimension” of a type system is well-defined.*

### 3.3 Compression Optimality

**Corollary 3.5** (Compression Optimality). *All minimal complete type systems achieve the same compression ratio. No type system can be strictly more efficient than another while remaining complete.*

This means: all observer strategies achieve the same compression ratio when using minimal complete axis sets. The difference between nominal-tag and interface-only observers lies in *witness cost*, not compression efficiency.

## 4 Kolmogorov Witness

### 4.1 Witness Cost for Type Identity

Recall from Section 2 that the witness cost  $W(P)$  is the minimum number of primitive queries required to compute property  $P$ . For type identity, we ask: what is the minimum number of queries to determine if two values have the same type?

**Theorem 4.1** (Nominal-Tag Observers Achieve Minimum Witness Cost). *Nominal-tag observers achieve the minimum witness cost for type identity:*

$$W(\text{type identity}) = O(1)$$

*Specifically, the witness is a single tag read: compare  $\text{tag}(v_1) = \text{tag}(v_2)$ .*

*Interface-only observers require  $W(\text{type identity}) = \Omega(n)$  where  $n$  is the number of interfaces.*

*Proof.* See Lean formalization: `proofs/nominal_resolution.lean`. The proof shows:

1. Nominal-tag access is a single primitive query

2. Interface-only observers must query  $n$  interfaces to distinguish all types
  3. No shorter witness exists for interface-only observers (by the information barrier)
- 

## 4.2 Witness Cost Comparison

Observer Class	Witness Procedure	Witness Cost $W$
Nominal-tag	Single tag read	$O(1)$
Interface-only	Query $n$ interfaces	$O(n)$

Table 1: Witness cost for type identity by observer class.

To our knowledge, this is the first formal proof that nominal-tag access minimizes witness cost for type identity.

## 5 Rate-Distortion Analysis

### 5.1 Three-Dimensional Tradeoff: Tag Length, Witness Cost, Distortion

Recall from Section 2 that observer strategies are characterized by three dimensions:

- **Tag length  $L$ :** machine words required to store a type identifier per value
- **Witness cost  $W$ :** minimum number of primitive queries to implement type identity checking
- **Distortion  $D$ :** worst-case semantic failure flag ( $D = 0$  or  $D = 1$ )

We compare two observer classes:

**Definition 5.1** (Interface-only observer). An observer that queries only interface membership ( $q_I \in \Phi_{\mathcal{I}}$ ), with no access to explicit type tags.

**Definition 5.2** (Nominal-tag observer). An observer that may read a single type identifier (nominal tag) per value, in addition to interface queries.

**Theorem 5.3** (Pareto Optimality of Nominal-Tag Observers). *Nominal-tag observers achieve the unique Pareto-optimal point in the  $(L, W, D)$  space:*

- **Tag length:**  $L = O(1)$  machine words per value
- **Witness cost:**  $W = O(1)$  primitive queries (one tag read)
- **Distortion:**  $D = 0$  (type equality implies behavior equivalence)

Interface-only observers achieve:

- **Tag length:**  $L = 0$  (no explicit tag)
- **Witness cost:**  $W = O(n)$  primitive queries (must query  $n$  interfaces)
- **Distortion:**  $D = 1$  (type equality does not imply behavior equivalence)

*Proof.* See Lean formalization: `proofs/python_instantiation.lean`. The proof verifies:

1. `nominal_cost_constant`: Nominal-tag achieves  $(L, W, D) = (O(1), O(1), 0)$
2. `duck_cost_linear`: Interface-only requires  $O(n)$  queries
3. `python_gap_unbounded`: The cost gap is unbounded in the limit
4. Interface observations alone cannot distinguish provenance; nominal tags can

■

## 5.2 Pareto Frontier

The three-dimensional frontier shows:

- Nominal-tag observers dominate interface-only observers on all three dimensions
- Interface-only observers trade tag length for distortion (zero  $L$ , but  $D = 1$ )

To our knowledge, this is the first formal proof of Pareto optimality for nominal-tag observers in the  $(L, W, D)$  tradeoff.

*Remark 5.4* (Programming language instantiations). In programming language terms: *nominal typing* corresponds to nominal-tag observers (e.g., CPython’s `isinstance`, Java’s `.getClass()`). *Duck typing* corresponds to interface-only observers (e.g., Python’s `hasattr`). *Structural typing* is an intermediate case with  $D = 0$  but  $W = O(n)$ .

# 6 Instantiations in Real Runtimes

## 6.1 CPython: Nominal-Tag Access

**Corollary 6.1** (CPython instantiation). *CPython realizes nominal-tag access: the runtime stores a type tag (`ob_type` pointer) per object and exposes it via the `type()` builtin. Therefore, the constant-witness construction applies:  $W(\text{type-identity}) = O(1)$  in CPython.*

**Evidence:** The CPython object model stores a pointer to the type object (`ob_type`) in every heap-allocated value [1]. The `type()` builtin [2] is a single pointer dereference, hence  $O(1)$  primitive queries.

## 6.2 Java: Nominal-Tag Access

**Corollary 6.2** (Java instantiation). *Java realizes nominal-tag access via the `.getClass()` method, which returns the runtime type object. Like CPython, Java achieves  $W(\text{type-identity}) = O(1)$ .*

## 6.3 TypeScript: Structural Typing

**Corollary 6.3** (TypeScript instantiation). *TypeScript uses structural typing: two types are equivalent iff they have the same structure (fields and method signatures). Type identity checking requires traversing the structure, hence  $W(\text{type-identity}) = O(n)$  where  $n$  is the structure size.*

## 6.4 Rust: Nominal-Tag Access (Compile-Time)

**Corollary 6.4** (Rust instantiation). *Rust uses nominal typing at compile time: type identity is resolved statically via the type system. At runtime, Rust provides `std::any::type_id()` for nominal-tag access, achieving  $W(\text{type-identity}) = O(1)$ .*

## 6.5 Summary: Witness Cost Across Runtimes

Language	Observer Class	Witness Cost $W$
Cython	Nominal-tag	$O(1)$
Java	Nominal-tag	$O(1)$
TypeScript	Interface-only (structural)	$O(n)$
Rust	Nominal-tag	$O(1)$

Table 2: Witness cost for type identity across programming language runtimes.

These instantiations confirm the theoretical predictions: nominal-tag observers achieve constant witness cost, while interface-only observers require linear witness cost in the number of interfaces.

## 7 Conclusion

This paper presents an information-theoretic analysis of semantic inference under observational constraints. We prove three main results:

1. **Impossibility Barrier:** No interface-only observer can compute properties that vary within indistinguishability classes.
2. **Constant-Witness Result:** Nominal-tag observers achieve  $W(\text{type-identity}) = O(1)$ , the minimum witness cost.
3. **Pareto Optimality:** Nominal-tag observers achieve the unique Pareto-optimal point in the  $(L, W, D)$  tradeoff: minimal tag length, minimal witness cost, zero distortion.

### 7.1 Implications

These results have several implications:

- **Nominal-tag observers are provably optimal** for type identity checking, not just a design choice.
- **Interface-only observers are provably limited:** they cannot achieve  $D = 0$  regardless of computational resources.
- **The barrier is informational, not computational:** even with unbounded time and memory, interface-only observers cannot overcome the indistinguishability barrier.

*Remark 7.1* (Programming language instantiations). In PL terms: nominal typing (Cython, Java) instantiates nominal-tag observers; duck typing (Python `hasattr`) instantiates interface-only observers; structural typing is intermediate ( $D = 0, W = O(n)$ ).

## 7.2 Future Work

This work opens several directions:

1. **Other observation families:** Do other semantic concepts (modules, inheritance, generics) induce matroid structure on their observation spaces?
2. **Witness complexity of other properties:** What are the witness costs for provenance, mutability, or ownership semantics?
3. **Hybrid observers:** Can we design observer strategies that achieve better  $(L, W, D)$  tradeoffs by combining tag and interface queries?

## 7.3 Conclusion

Semantic inference under observational constraints admits a clean information-theoretic analysis. Nominal-tag observers are not merely a design choice—they are the provably optimal strategy for type identity under the  $(L, W, D)$  tradeoff. All proofs are machine-verified in Lean 4.

## References

- [1] Python Software Foundation. Cpython object implementation: Object structure. <https://github.com/python/cpython/blob/main/Include/object.h>, 2024.
- [2] Python Software Foundation. Python data model: type(). <https://docs.python.org/3/reference/datamodel.html>, 2024.

## A Lean 4 Formalization

All theorems in this paper are formalized and machine-verified in Lean 4. The proofs are located in the repository at:

`docs/papers/paper1_typing_discipline/proofs/`

### A.1 Proof Statistics

- **Total Lines:** ~6,000
- **Theorems:** 265
- **Lemmas:** 150+
- **Sorry Placeholders:** 0 (all proofs complete)
- **Axioms Used:** propext (proposition extensionality)

### A.2 Key Proof Files

1. `abstract_class_system.lean`: Core formalization of the class system model, interface equivalence, and information barrier theorem
2. `axis_framework.lean`: Type axis matroid structure, equicardinality proofs (`semantically_minimal_implies_orthogonal`, `minimal_complete_unique_orthogonal`)

```
3. python_instantiation.lean:      Witness      cost      proofs      (nominal_cost_constant,
   duck_cost_linear, python_gap_unbounded)
```

*Remark A.1.* The key theorems referenced in this paper are distributed across these files. The paper cites specific lemma names to enable direct verification.

### A.3 Building the Proofs

To verify the proofs locally:

```
cd docs/papers/paper1_typing_discipline/proofs
lake update
lake build
```

All theorems will be machine-verified if compilation succeeds with no errors.

### A.4 Axiom Dependencies

The proofs use only one axiom: `propext` (proposition extensionality). This is a standard axiom in constructive mathematics and does not affect the validity of the results.

All other proofs are constructive (no use of `Classical.choice` or `Decidable.em`).

### A.5 Reproducibility

The Lean toolchain version is specified in `lean-toolchain`. All dependencies are pinned in `lake-manifest.json`. The proofs are reproducible on any system with Lean 4 installed.