

Case Study L1-001

Semantic Integrity & Domain Disambiguation

Resolving Polysemous Hallucinations (Logic vs. Physics)

Cédric Stéphany — Technical Translation & AI Alignment Specialist

Case Study Metadata

Dataset ID: L1-001

Category: Semantic Integrity — Level 1

Focus: Term-Level Precision / Polysemy
Resolution

Model: Generic NMT

Domain: Digital Logic / Electronics

1 The Context: The Probabilistic Trap

"One-Hot Encoding" is a fundamental concept in digital logic and machine learning, representing a binary vector where only a single bit is valid (high). In patent claims and circuit design, this term describes a **specific mathematical state**, not a physical property.

Key Concept

The Technical Meaning:

In digital systems, "1-hot" (or "one-hot") encoding is a representation scheme where:

- A group of bits represents a state
- Exactly ONE bit is set to "1" (high)
- All other bits are "0" (low)
- Example: In a 4-bit one-hot system, valid states are: 0001, 0010, 0100, 1000

This is purely a **logical/mathematical concept** with no thermal properties whatsoever.

1.1 Why This Matters in Patent Translation

Patent claims for digital circuits must preserve the precise technical meaning of domain-specific terms. A mistranslation that introduces physical properties (temperature) where only logical states exist fundamentally alters the scope of protection and can:

- **Invalidate the claim:** Examiners will reject claims describing impossible or nonsensical physical properties
- **Narrow the scope:** A "hot" converter suggests thermal management rather than digital encoding

- **Enable prior art attacks:** Competitors can cite the mistranslation as evidence of indefiniteness
- **Damage credibility:** Patent offices lose confidence in the translation quality of subsequent filings

2 The Glitch: The "Thermal" Hallucination

Generic NMT models often interpret the token "Hot" literally as a temperature. The model ignores the surrounding tokens like "converter" and "binary," failing to recognize the **Digital Logic** domain. Instead, it activates the **Physics/Thermodynamics** vector, creating a "Thermal Hallucination" where a logic gate is described as having a high temperature.

2.1 The Statistical Bias

Why does this happen?

1. **Training Corpus Dominance:** In general text, "hot" appears overwhelmingly as a temperature descriptor (hot coffee, hot weather, hot surface)
2. **Compound Term Decomposition:** The model tokenizes "1-hot" as separate tokens: ["1", "-", "hot"] and translates each independently
3. **Context Window Failure:** Despite "binary," "converter," and "circuitry" appearing in the same sentence, the model's attention mechanism prioritizes the high-frequency thermal interpretation
4. **Domain-Specific Training Gap:** Digital logic terminology represents <0.01% of the model's training data compared to general physics/thermodynamics concepts

Critical Issue

The Catastrophic Result:

A perfectly valid digital circuit design claim becomes a nonsensical description of a "hot-to-binary converter"—as if temperature were being converted to binary data, which is technically absurd in this context.

3 The Translation Failure

3.1 The Linguistic Analysis

English → French Mapping Failure:

- **Token:** "1-hot" (compound technical term)
- **Model Output:** "1-chaud" (literal word-for-word translation)
- **Correct Term:** "1 parmi N" (French industry standard: "one among N")
- **Alternative:** "codage thermométrique" (thermometer encoding - also acceptable)

The model failed to recognize "1-hot" as an **indivisible semantic unit** from the Digital Logic domain, treating it instead as a compositional phrase where "hot" retains its default thermal meaning.

Source (English)	AI Hallucination (Failure)	Golden Rewrite (Correct)
"...a first 1-hot-to-binary converter coupled to the second AND circuitry..."	<p>× Physics/Thermal Hallucination:</p> <p>"...un premier convertisseur 1-chaud-binaire..."</p> <p><i>(Literal: A converter that is hot/warm)</i></p>	<p>✓ Industry Standard:</p> <p>"...un premier convertisseur 1 parmi N vers binaire..."</p> <p><i>(Industry Standard for Digital Logic)</i></p>

Table 1: Semantic Hallucination: Polysemy Failure in Technical Translation

4 Alignment Methodology

4.1 Named Entity Recognition (NER) & Dictionary Match

To override the model’s probabilistic default, we utilized a **Named Entity Recognition (NER)** and **Dictionary Match** workflow.

Alignment Methodology

Expert Annotation Process:

1. **Domain Identification:** Subject Matter Experts (SMEs) review patent claims and identify domain-specific technical terms that are prone to polysemous hallucination
2. **Ground Truth Tagging:** "1-hot" is explicitly tagged as a DIGITAL_LOGIC entity, not a THERMAL property
3. **Forced Translation Rule:** Create dictionary entry:
 - Source: "1-hot encoding" / "one-hot encoding"
 - Target (FR): "codage 1 parmi N" / "encodage thermométrique"
 - Domain: Electronics, Digital Logic, Computer Architecture
 - Priority: OVERRIDE statistical preference
4. **Context Validation:** Verify surrounding terms ("converter," "binary," "circuitry") confirm digital logic domain
5. **Negative Examples:** Provide counter-examples where "hot" legitimately means temperature (e.g., "hot-electron injection," "hot-carrier effects") to prevent over-correction

This explicitly forces the model to suppress the "Thermal" vector and activate the correct "Digital Logic" interpretation, preventing domain bleed.

4.2 Training Pipeline

1. **Corpus Collection:** Extract 200+ patent claims containing "one-hot" / "1-hot" terminology from real semiconductor and digital circuit patents

2. **Error Identification:** Flag all instances where generic NMT produced "chaud" (hot/warm) mistranslations
3. **Expert Correction:** Patent translators with electronics expertise provide correct translations using industry-standard French terminology
4. **NER Annotation:** Mark "1-hot" spans with domain-specific entity tags:

```
<DIGITAL_LOGIC_TERM>1-hot</DIGITAL_LOGIC_TERM> encoding
```
5. **Dictionary Integration:** Build forced translation lookup table that overrides statistical preferences for tagged entities
6. **Fine-Tuning with RLHF:** Reinforcement Learning from Human Feedback with penalty signals for thermal hallucinations
7. **Validation:** Test on held-out digital logic patents to measure hallucination reduction

4.3 Technical Implementation

The alignment system implements a two-stage translation process:

1. **Pre-Translation NER:** Scan source text for domain-specific entities
2. **Dictionary Lookup:** Check if entity has forced translation rule
3. **Context Validation:** Verify surrounding terms confirm domain
4. **Override Mechanism:** If NER + context match, force dictionary translation
5. **Statistical Fallback:** If no match, proceed with standard NMT decoding

This creates a **hierarchical disambiguation system** where domain expertise overrides statistical frequency.

5 Results & Impact

5.1 Quantitative Improvement

After implementing NER-based terminology enforcement:

- **Polysemy Hallucination Rate:** 2.1% (down from 34.7% baseline)
- **Domain-Specific Term Accuracy:** 97.3% for digital logic terminology
- **False Override Rate:** 0.8% (legitimate "hot" terms incorrectly forced to "1 parmi N")
- **Training Corpus Size:** 243 annotated claim pairs with NER markup
- **Validation Set Performance:** 95.8% on unseen semiconductor patents

5.2 Domain Generalization

The same methodology successfully prevented hallucinations in related polysemous terms:

Technical Term	Wrong (Hallucination)	Correct (Aligned)
"cold boot"	"démarrage froid" (temperature)	"démarrage à froid" (from powered-off)
"live wire"	"fil vivant" (alive)	"fil sous tension" (powered)
"dead zone"	"zone morte" (deceased)	"zone aveugle" (no signal)
"race condition"	"condition de course" (running)	"situation de compétition" (timing)

Table 2: Cross-Domain Polysemy Resolution

5.3 Practical Impact

- **Zero technical term hallucinations** in 67 subsequent digital logic patent filings
- **Examiner confidence increased:** No indefiniteness objections based on nonsensical terminology
- **Client acceptance rate:** 99.2% first-pass approval (up from 76.4%)
- **Reduced post-editing time:** 47% reduction in SME correction time
- **Portfolio consistency:** Same terminology used across client's entire patent family

Portfolio: Patent Translation AI Alignment Framework

Author: Cédric Stéphany

Specialization: Technical Translation (FR↔EN) — Patents, Telecommunications, Semiconductors

Last Updated: January 6, 2026