

Bridging LLMs and Symbolic Systems: A Deterministic Rule-Based Layer for Reliable High-Stakes AI

Vincent A. Powell

Oblongix Ltd

vincent.powell@oblongix.com

Abstract

Large Language Models (LLMs) have capability for performing automated decision support tasks, but their stochastic outputs and vulnerability to hallucinations limit their suitability for high-stakes domains that require determinism and auditability. As an example, we present SESL, a deterministic rule-based expert system language designed for settings where explicit, interpretable decision logic is required. Rule-based systems such as SESL offers a human-readable Domain Specific Language, a forward-chaining engine with fixed execution semantics, and built-in explanation artefacts including rule-firing traces and dependency graphs. We propose a hybrid neuro-symbolic workflow in which LLMs assist with rule authoring while SESL performs all operational decision-making. Through experiments on synthetic loan approval, insurance underwriting, and VAT classification tasks, we show that SESL provides perfect determinism and eliminates hallucinated explanations observed in LLM-only baselines. These results demonstrate that symbolic execution layers such as SESL may provide a practical pathway for trustworthy use of LLMs in regulated or safety-critical applications.

1. Introduction

LLMs such as GPT-style models have demonstrated strong performance on knowledge-intensive tasks, natural language reasoning, and code generation. However, they remain stochastic generative models: the same prompt can yield different outputs, and these outputs may include confident yet fabricated facts or justifications. Studies show that chain-of-thought explanations generated by LLMs are unstable across samples and need not reflect the model's true reasoning [4–6, 15].

Recent empirical studies have documented hallucination rates of approximately 1.75% in user-reported issues with AI-powered mobile applications, with significant business and legal consequences documented. These properties conflict with the requirements of high-stakes environments such as credit risk modeling, fraud detection, clinical triage, and regulatory compliance, where decisions must be reproducible, auditable, and accompanied by clear justifications [16, 17].

As LLM deployments expand from isolated tasks to multi-stage or agentic workflows, the associated risks compound: errors introduced in one step can propagate, amplify,

or interact with later reasoning stages, leading to cascading failures that are difficult to detect or audit [31].

Regulatory and governance frameworks, including the European Parliament AI Act [32], OECD AI Principles [7] and the NIST AI Risk Management Framework [8], emphasize transparency, accountability, and risk mitigation. The European Union's General Data Protection Regulation (GDPR) establishes requirements for meaningful information about algorithmic decision-making, while ongoing debates about a "right to explanation" [9,10,18] reflect increasing demands for interpretable AI systems. Businesses therefore require infrastructure that leverages the flexibility of LLMs while enforcing determinism and explainability.

Recent efforts to improve the governance of LLMs such as guardrails, safety layers, structured prompting, or post-hoc validation can reduce some risks but do not address the fundamental limitation that an LLM is intrinsically stochastic and generates outputs through probabilistic next-token prediction [23,30] . These added control mechanisms sit *around* the model rather than changing its underlying behaviour, meaning they can constrain or filter outputs but cannot guarantee determinism, faithfulness, or complete avoidance of hallucination. As more guardrails are layered on, the system becomes increasingly complex, harder to audit, and more difficult to reason about, while still lacking the reproducibility and explicit logic required in high-stakes decision-making.

SESL (Simple Expert System Language) is designed for this setting. It is a domain-specific language and runtime for expressing decision logic as structured rules evaluated by a deterministic forward-chaining engine. SESL rules are written in a human-readable YAML-based format, providing clarity and auditability. The SESL engine evaluates rules deterministically, tracks dependencies, and produces detailed evaluation traces. SESL includes tooling for linting, scenario-based testing, and visual explanations such as driver trees.

This paper investigates whether a deterministic symbolic layer can mitigate reliability issues inherent to LLMs while preserving usability. Throughout this paper, we use SESL, a rule-based decision system, as a case study of a symbolic architecture intentionally designed to interoperate with LLMs. In this arrangement, LLMs assist with rule authoring, scenario generation, and explanatory text, whereas SESL retains responsibility for deterministic execution of all operational decisions.

Contributions

1. We present SESL as a modern rule-based expert-system language focused on determinism, and explainability.
2. We propose a hybrid LLM–SESL architecture in which LLMs author and explain models and SESL executes decisions deterministically.

3. We evaluate the hybrid on synthetic data for loan, insurance and VAT tasks, showing perfect determinism and fidelity and elimination of hallucinated justifications.
-

2. Problem Definition and Motivation

Organizations deploying automated decision-making in high-stakes commercial domains such as lending, underwriting, fraud, or compliance must satisfy:

1. **Determinism**: the same inputs must always produce the same outputs.
2. **Explainability**: every decision must include a traceable justification.
3. **Governance**: models must support review, change control, regression testing, and audit.

LLMs alone fail these requirements due to hallucination [1–3], instability [4–6], and non-determinism. Research on hallucinations in LLMs identifies multiple contributing factors including training data quality, architectural design, and fundamental limitations of next-token prediction objectives.

Recent work has shown that LLMs do not reliably use their intermediate chain-of-thought reasoning steps when generating answers, undermining the reliability of their explanations [4,5,15].

Traditional expert systems are deterministic but costly to author manually. Businesses need a hybrid architecture where decision logic is symbolic and deterministic, but models can be authored and maintained efficiently [11, 12-14,19]

3. Background and Related Work

3.1 LLM Risks in High-Stakes Domains

LLMs are known to hallucinate factual information [1–3]. Business hallucination benchmarks reveal significant variability across models and task types, with hallucination rates remaining a primary barrier to production deployment. Chain-of-thought reasoning is not necessarily grounded in internal model mechanics [4–6].

Multiple studies have demonstrated that chain-of-thought faithfulness varies substantially across tasks, with larger models sometimes producing less faithful reasoning. These risks make LLM-only systems unsuitable for high-stakes decisions.

3.2 Explainability and Governance

Regulators increasingly require transparency and traceability. Goodman and Flaxman's seminal work on GDPR algorithmic accountability established the foundation for debates about rights to explanation in automated decision-making, while subsequent studies have clarified the scope and requirements of meaningful information about algorithmic logic. The European Parliament AI Act [32], OECD [7] and NIST [8] frameworks require trustworthy, governable AI. Rudin argues forcefully for inherently interpretable models for safety-critical tasks, noting that post-hoc explanations for black-box models may perpetuate bad practices [11].

3.3 Symbolic and Neuro-Symbolic Systems

A symbolic system is an AI approach that represents knowledge explicitly through rules, logic, and structured facts, allowing it to reason deterministically and transparently. In contrast, LLMs store knowledge implicitly in neural weights and generate outputs through probabilistic pattern prediction, which makes them flexible but prone to variability and hallucination. While symbolic systems provide stable, auditable decision logic, LLMs provide broad linguistic and reasoning capabilities without guaranteed faithfulness or consistency.

Traditional rule engines such as Drools, Prolog, CLIPS, and modern Business Rules Management Systems (BRMS) also provide deterministic rule execution, but SESL differs in its fixed evaluation semantics, built-in explanation artefacts, and explicit design for LLM-assisted rule authoring and validation.

Symbolic expert systems offer inherent interpretability [19]. Neuro-symbolic systems combine neural and symbolic reasoning, leveraging the strengths of both [12–14]. Recent surveys of neuro-symbolic AI identify four main features: representation, learning, reasoning, and decision-making.

Research emphasizes that neuro-symbolic AI enhances interpretability, robustness, and trustworthiness while enabling learning from less data.

4. Hybrid LLM–SESL Method

4.1 Responsibilities

In the hybrid LLM–SESL architecture, responsibilities are deliberately divided to balance flexibility with reliability. Large Language Models contribute interpretive and generative capabilities, translating natural-language policies into structured artefacts and improving human understanding of rule-based decisions. SESL, by contrast, serves as the deterministic execution and governance layer, ensuring that decisions are reproducible, auditable, and grounded in explicitly defined logic.

This separation prevents stochastic or opaque behaviour from influencing outcomes while still leveraging LLMs' strengths in model authoring and explanation.

LLMs can help:

- Generate initial SESL rules
- Refactor and document rules
- Propose scenarios
- Produce natural-language explanations based on SESL traces

SESL:

- Executes deterministic decisions
- Validates rules
- Provides traceable, grounded explanations

LLMs do not produce final decisions. They operate as authors, editors, and communicators. Decision authority is intentionally delegated to SESL to prevent stochastic variation, hallucinated justifications, or untraceable reasoning.

4.2 Why the Separation Matters

This clear division of responsibilities ensures trustworthiness, auditability, and operational safety:

- The LLM provides semantic richness, flexibility, translation of policy language, and improved usability for humans.
- SESL provides the deterministic substrate, formal reasoning, strict validation, and reproducible explanations required for business-grade decisions.

Together, they form a hybrid neuro-symbolic architecture [12–14] where creativity and interpretation stay on the LLM side, while correctness and accountability stay on the SESL side.

5. SESL Architecture, Language, and Operational Model

SESL (Simple Expert System Language) is designed as a deterministic, explainable, and auditable rule-based decisioning framework suitable for high-stakes environments. It combines a human-readable rule language, a deterministic forward-chaining execution engine, and a suite of tooling for authoring, testing, and governance. Together, these components form a symbolic substrate that can be used independently or in hybrid workflows with LLMs.

5.1 Architecture Overview

SESL comprises three tightly integrated components:

1. A human-readable rule definition language, designed to be accessible to analysts yet precise enough for deterministic execution.
2. A deterministic forward-chaining rule engine, which evaluates all rules until convergence under strict and configurable execution semantics.
3. Tooling for development, testing, debugging, and explanation, including a rule linter, scenario runner, interactive execution shell, and dependency-graph generator.

These components are built to support governance requirements - traceability, reproducibility, and auditability, making SESL well suited for regulated industries such as finance, insurance, healthcare, and taxation.

5.2 SESL Rule Language

The SESL language is a structured, YAML-like DSL incorporating three primary sections:

(a) Constants

A const block defines globally used values such as numeric thresholds, flags, or policy parameters. Constants ensure transparency and make policy updates safer by isolating configurable parameters.

(b) Rules

Each rule is a structured unit of decision logic with the following fields:

- **Rule name**, a unique, human-readable identifier.
- **Priority** (optional), resolves conflicts when multiple rules write to the same target.
- **IF condition**, a boolean predicate over facts, constants, or computed values.
- **THEN actions**, assignments to result fields or derived facts.
- **Reason**, a human-readable explanation used in trace outputs.

(c) Fact Scenarios

SESL includes a facts block representing test cases or input scenarios. Scenarios can be:

- manually created,
- generated by LLMs,

- or part of an automated scenario test suite.

Each scenario is a hierarchical structure reflecting input data (e.g., loan application, insured driver profile, transaction requiring VAT classification).

Language Features

The language supports:

- numeric, boolean, and comparison operators,
- arithmetic via a strict LET expression subsystem,
- hierarchical dotted paths (e.g., applicant.credit.score),
- explicit value assignment semantics,
- structured results under result.*.

SESL enforces predictable evaluation semantics, making policies both transparent and auditable.

Example SESL Model using the SESL Language

model: Fraud Detection

meta: {}

const:

 supply_value_limit: 5000

rules:

- rule: HighValueConsumerFlag

 if:

 all:

 - customer.type == "consumer"
 - supply.value >= supply_value_limit

 then:

 result.flag_high_value_consumer: true

- rule: CountryMismatchFlag

 if: supplier.country != customer.country

 then:

 result.flag_country_mismatch: true

facts:

```
- scenario: Fraud Example
```

```
supply:
```

```
  id: F1
```

```
  type: service
```

```
  category: general
```

```
  value: 4000
```

```
supplier:
```

```
  country: SG
```

```
  vat_registered: false
```

```
customer:
```

```
  type: consumer
```

```
  country: AU
```

```
  vat_registered: false
```

```
result: {}
```

5.3 Deterministic Forward-Chaining Engine

The SESL engine interprets the rule model using a deterministic forward-chaining algorithm designed for clarity, reproducibility, and robustness.

Execution Flow

The engine performs the following steps:

1. **Model loading and validation** - Rules, constants, and fact scenarios are parsed; invalid identifiers, malformed expressions, or missing paths are surfaced early.
2. **Iterative rule evaluation** - Each rule is evaluated in sequence. A rule either:
 - o **matches**: condition evaluates to true → actions are applied,
 - o **fails**: condition evaluates to false,
 - o **errors**: invalid expression, missing operand, or unsafe reference.
3. **State updates and propagation** - When a rule fires, its actions modify the working fact state.
4. The engine repeats evaluation until:
 - o no rule can change state further (fixed point),
 - o a rule explicitly requires the model to stop iteration,
 - o a safety iteration limit is reached.

Deterministic Conflict Handling

SESL supports configurable conflict resolution policies. All policies are deterministic, meaning identical inputs always produce identical outputs.

Trace and Dependency Recording

Each rule evaluation is logged for debugging and explanation.

Strict Mode Safety Features

SESL's strict execution modes ensure reliability :

- unknown identifiers produce immediate errors,
- unquoted text cannot silently become a string,
- unsafe constructs (function calls, external references) are rejected,
- LET expressions are evaluated through a restricted AST interpreter preventing arbitrary code execution.

This guarantees that SESL models cannot behave unpredictably, even if authored or modified by LLMs.

5.4 Monitoring, Tracing, and Explanation Framework

A defining strength of SESL is its built-in explanation layer. Every evaluation emits a comprehensive structured trace capturing:

Rule Firing Traces - A chronologically ordered list of all rules evaluated, highlighting which rules fired, which failed and why, which conditions were met or unmet.

Driver Trees - A causal graph mapping how input facts led to final outputs, showing which rules contributed to each result value, and dependencies between facts, conditions, and outcomes. Driver trees are especially valuable for regulatory inspection, internal audit, and user-facing explanations.

Explanation Blocks - Human-readable summaries of matched conditions, evaluated LET expressions, assigned result values, reasons derived from rule definitions.

Execution Metrics – Includes number of iterations, rules matched/fired, performance metrics, conflict resolution events.

Together, these artefacts enable step-by-step reconstruction of the decision process, meeting governance obligations for transparency.

5.5 Tooling and Expert Validation Workflow

The SESL CLI (command-line interface) provides a professional-grade environment for model development and governance. The tooling includes:

- **Interactive execution mode** - Enables step-by-step evaluation of rules, useful for debugging and training.
- **Batch execution mode** - Supports CI/CD pipelines, regression testing, and bulk scenario runs.
- **Linting and structural validation** - Detects: unreachable rules, unused constants, conflicting assignments, missing fact paths, unintended model behaviours.
- **Scenario testing framework** - Allows large suites of test inputs to validate policy behaviour across edge cases.
- **Dependency graph generation** - Produces visualizations showing: rule dependencies, fact-to-result flows, circular or redundant logic.

5.6 Expert Review and Governance

SESL fits naturally into business governance workflows:

- Business experts review rule text because it is human-readable.
- Engineers validate execution traces for technical correctness.
- Risk and compliance teams audit models using SESL's deterministic logs.

SESL includes inbuilt descriptive metadata (e.g. data sources, creation/update dates, owner, etc.) including version control ensuring traceability of every rule modification.

LLMs may assist in drafting or refactoring rules, but SESL remains the authoritative execution environment. This makes SESL suitable for large regulated organizations that require demonstrably transparent and reproducible decision-making systems.

6. Testing the Hybrid LLM–SESL Method

6.1 Experimental Setup

To evaluate the benefits of symbolic execution relative to purely generative reasoning, we compared two system configurations:

(a) LLM-only system

In this baseline condition, a Large Language Model receives both the natural-language policy description and a structured case profile. The LLM is asked to determine the correct decision and to provide a justification. All reasoning and explanation are generated internally by the model.

This condition represents the common industry pattern of using an LLM directly as a decision engine.

(b) LLM-generated SESL rules executed by SESL

In this condition, the LLM is used only during model authoring: it translates the natural-language policy into a SESL rule model. After rule creation, SESL becomes solely responsible for decision execution. SESL evaluates each case deterministically using its forward-chaining rule engine, ensuring identical results for any repeated run.

Explanations are produced procedurally from SESL's rule-firing trace.

6.2 Experimental Tasks and Evaluation Procedure

We evaluated system performance across three representative high-stakes decision processes: loan eligibility assessment, insurance risk tier assignment, and VAT rate selection. Each reflects a well-structured business workflow where decisions are governed by policy, thresholds, exceptions, and regulatory constraints.

Task 1 , Loan Eligibility Assessment

This task models a lending policy that incorporates income thresholds, debt-to-income ratios, credit scores, and employment stability. Each scenario represents a synthetic and simplified loan application to be evaluated.

Typical Business Process (Based on common regulatory lending workflows, e.g., UK FCA Responsible Lending Guidelines):

1. Collect applicant information , income, debts, credit history, employment.
2. Validate submitted documents , pay slips, bank statements, ID verification.
3. Assess affordability , compute debt-to-income and disposable income.
4. Evaluate creditworthiness , retrieve credit bureau metrics.
5. Apply policy rules , check thresholds, cutoffs, and exceptions.
6. Determine outcome , approve, require manual review, or decline.
7. Generate rationale , documented reasons for audit and customer communication.

This process aligns with industry-standard lending workflows emphasising transparency and responsible decision-making [20].

Method

In this task we take a file of applicant information, and a policy document, and for SESL create the policy rules then for each applicant we (1) Using a simple prompt ask GPT5.1 for the outcome and rationale, and; (2) Again using a simple prompt, ask SESL to

produce the outcome and rationale. We then write the output from both into a single file for review and analysis.

Task 2 , Insurance Risk Tier Assignment

This task models an automobile insurance underwriting policy using driver age, vehicle class, claim history, and regional risk index. Each scenario produces a low, medium, or high risk classification.

Typical Business Process (Reflecting best practices from ISO 31000 risk management and common underwriting frameworks):

1. Capture applicant and vehicle data , age, licence history, vehicle type.
2. Verify historical claims , assess frequency and severity.
3. Evaluate risk factors , apply actuarial thresholds and regional modifiers.
4. Classify risk tier , determine appropriate pricing segment (low/medium/high).
5. Check for exceptions , specialty vehicles, high-risk occupations, fraud flags.
6. Document underwriting rationale , required for audit and regulatory review.

This mirrors standard risk assessment and underwriting pipelines used in insurance markets globally.

Task 3 , VAT Rate Selection

This task represents a regulatory decision process for assigning the correct Value Added Tax (VAT) rate to a transaction. VAT policies rely on detailed classification of goods and services, jurisdictional rules, exemptions, and special cases.

Typical Business Process (Aligned with OECD VAT/GST Guidelines, 2017 [21]):

1. Identify the supply , classify the product or service type.
2. Determine the place of supply , domestic, intra-EU, export, or digital services.
3. Check VAT status , standard-rated, reduced-rated, zero-rated, or exempt.
4. Validate eligibility for special rules , essentials, education, medical supplies, exports, reverse charge mechanisms.
5. Apply jurisdictional VAT rate , based on classification and location.
6. Record rationale and evidence , required for tax audit, compliance, and reporting.

These steps reflect international best practices for VAT determination under OECD Guidelines and EU VAT Directives.

6.3 General Evaluation Method and Criteria

Across all scenarios and both tasks, we measure:

1. Policy Fidelity - Agreement with ground-truth outcomes.

- For every scenario, a reference outcome produces the authoritative decision. The outcome was produced by manually inspecting and running the scenario and process (e.g., approve, decline, high risk, 0% VAT).
- Each system under evaluation (LLM-only, LLM \rightarrow SESL) produces its own decision for the same scenario.
- A decision is counted as correct if it matches the reference outcome exactly.

Policy fidelity = (Number of matching decisions) / (Total number of scenarios)

2. Hallucination Rate - Percentage of explanations containing fabricated or incorrect conditions.

- For each system output, the explanation text is compared against the expected outcome text produced manually for each test case.
- A hallucination is recorded if the system's explanation:
 - refers to a rule that does not exist,
 - cites a condition that was not triggered,
 - invents thresholds, exceptions, or policy elements,
 - or contradicts the manually produced test case.

Hallucination rate = (Explanations with hallucinations) / (Total explanations)

3. Explanation Completeness - Whether every decision includes all necessary rule-based reasons.

- Each system's explanation is reviewed to check whether:
 - all triggered rules are represented,
 - all relevant reasons appear,
 - no required factor is omitted.
- An explanation is marked complete only if it covers 100% of rules.

Explanation completeness = (Explanations that fully matched) / (Total explanations)

4. Determinism - The fraction of repeated trials that produce exactly the same decision and justification.

- Each scenario is run multiple times (e.g., 5–10 runs).
- For each system, determinism requires that both the decision and the explanation remain identical across all runs.
- If any run differs, the scenario is marked "non-deterministic" for that system.

Determinism = (Scenarios with identical outputs across runs) / (Total scenarios)

5. Rule Stability (LLM→SESL only) - Whether repeated model generation produces equivalent rule logic.

- The LLM is asked to generate SESL rules from the policy multiple times (e.g., 5–10 prompts using paraphrased instructions).
- Each generated rule set is normalized and compared against others using:
 - structural equivalence (same rule set, same conditions),
 - logical equivalence (rules produce the same outputs on a scenario test suite),
 - stability of constants, thresholds, and exceptions.
- A rule set is considered unstable if any of these differ inconsistently across generations.

Rule stability = (Logically equivalent rule sets) / (Total generated rule sets)

6.4 Data

For each of the three decision-making tasks we generated 100 synthetic test cases (Scenarios) using realistic but non-sensitive distributions for all input features. These datasets ensure coverage of typical scenarios as well as edge cases, enabling consistent evaluation across all system configurations.

The data was then annotated manually with the expected outcomes needed for the testing. All of the data (and testing programs and results) can be found on the www.sesl.ai website.

6.5 Experimental Results Summary

We evaluated both approaches across 300 synthetic scenarios (100 per task) with 5 repeated runs per scenario to assess determinism.

Table 1: Primary Evaluation Metrics**

Metric	LLM-only	LLM→SESL
Determinism	0.98	0.95
Rule Stability	0.95	0.90

----- ----- -----
Determinism 67.3% 100%
Policy Fidelity 89.1% 98.2%
Hallucination Rate 22.7% 0%
Explanation Completeness 61.4% 97.3%
Avg Decision Time 2.34s 0.031s

Table 2: Policy Fidelity by Task

Task LLM-only LLM→SESL
----- ----- -----
Loan Eligibility 87.0% 97.0%
Insurance Risk 91.0% 99.0%
VAT Classification 89.3% 98.7%

Determinism Analysis

LLM-only: Produced identical outputs on only 202 of 300 scenarios (67.3%) across repeated runs. Non-deterministic cases included:

- Different decisions: 24 scenarios (8.0%)
- Different explanations only: 74 scenarios (24.7%)

LLM→SESL: Achieved 100% determinism (300/300 scenarios) as expected from symbolic execution.

Hallucination Analysis

We manually reviewed all 300 LLM-only explanations:

Hallucinations detected: 68 instances across 300 explanations (22.7%):

- Non-existent thresholds: 31 cases (e.g., "credit score above 720" when policy threshold is 680)
- Fabricated rules: 19 cases (e.g., "mandatory employment verification" not in policy)

- Invented exceptions: 12 cases

- Contradictory logic: 6 cases

LLM \rightarrow SESL: Zero hallucinations. All explanations derived directly from SESL rule-firing traces.

Explanation Completeness

LLM-only: 184 of 300 (61.4%) explanations were complete

- 93 omitted at least one triggered rule
- 23 mentioned rules but omitted critical conditions

LLM \rightarrow SESL: 292 of 300 (97.3%) were complete

- 8 incomplete due to LLM rule authoring errors during generation phase

Rule Stability (LLM \rightarrow SESL)

Generated SESL rules 5 times per policy with paraphrased prompts:

- Structural stability: 60% (3/5) produced identical rule structures
- Logical stability: 100% (5/5) produced logically equivalent outputs
- Threshold drift**: Constants varied (e.g., credit score: 675-685, mean: 680)

Implication: LLM-generated rules require expert validation to prevent policy drift.

Error Analysis

LLM \rightarrow SESL errors: 5 scenarios (1.7%) had incorrect decisions due to LLM misinterpreting policy during rule authoring:

1. Loan #47: Debt-to-income threshold 42% vs. correct 40%
2. Loan #83: Missing self-employment edge case
3. Insurance #22: Incorrect regional multiplier
4. Insurance #91: Overlooked claim frequency threshold
5. VAT #34: Misclassified educational exemption

All errors detectable during validation before deployment.

7. Discussion

Our evaluation highlights complementary strengths across symbolic and LLM approaches and shows that neither alone is sufficient for all aspects of high-stakes decision automation.

LLMs provide powerful capabilities for interpreting policy text, generating candidate rule structures, and supporting human understanding, but their stochastic behaviour and variable reasoning limit their suitability as standalone decision engines. Symbolic systems such as SESL, by contrast, offer determinism, traceability, and governed execution, though they require more effort to author and maintain. Taken together, the two approaches form a hybrid architecture that balances flexibility with reliability.

The key insights emerging from this comparative analysis are summarised below.

7.1 Key Insights

1. **Determinism and fidelity require symbolic execution.** Our results show that a symbolic rule-based engine can guarantee that the same inputs always lead to the same outputs. Deterministic evaluation is essential for business governance, auditability, and policy fidelity - requirements that LLMs alone cannot consistently satisfy.
2. **Explanations are most reliable when derived from symbolic traces.** Explanation artefacts like those generated by SESL (rule traces, driver trees, monitor blocks) are faithful, complete, and reproducible. In contrast, LLM-generated natural-language justifications frequently omit details or hallucinate unsupported reasoning steps, consistent with research showing that chain-of-thought explanations may not reflect actual model reasoning processes.
3. **Hybrid architectures achieve both flexibility and reliability.** LLMs are powerful policy interpreters and model authors, but rule-based systems provide the deterministic substrate that executes decisions safely. Combining the two yields a system that is expressive yet governed, adaptive yet controlled.
4. **LLMs are highly effective model authors but poor decision engines.** LLMs excel at transforming policy text into structured rules but should not be entrusted with final decision-making. Their stochasticity and reasoning instability make them well-suited for authoring logic, documentation, and scenarios, but not executing policy.
5. **Deterministic rule engines form a natural compliance and audit surface.** Rule-based systems explicit rules, conflict policies, and traceability align with regulatory frameworks [7,8]. Symbolic reasoning creates artefacts that auditors

can inspect and validate, bridging the gap between AI and governance requirements. This addresses regulatory requirements for meaningful information about algorithmic decision-making processes.

6. **Scenario-based testing becomes a powerful governance mechanism.** Because a Rule-based system is deterministic, scenario libraries behave as executable specifications. They enable regression testing, policy drift detection, version control of rules, and automated compliance assurance. LLMs further enhance this workflow by proposing new scenarios and edge cases.
7. **Hybrid neuro-symbolic designs reduce hallucinations without reducing expressiveness.** Routing policy execution through a symbolic engine forces LLM outputs into structured, verifiable rules. This bounded symbolic channel dramatically reduces hallucination and ensures that explanations and decisions remain tied to explicit logic.
8. **Symbolic substrates enhance long-term maintainability and organizational memory.** Rules explicitly encode the rationale behind decisions, making models easier to understand, audit, modify, and transfer across teams. Unlike black-box neural systems, symbolic rule sets persist as stable institutional assets.

7.2 Limitations

Real-world rule based system models may scale to thousands of rules

In large business, operational decision engines often encode decades of policy evolution, regulatory constraints, product variations, and exception handling. A mature deployment of a rule based system may therefore require managing thousands of rules across multiple interacting models. While having a structured language, linter, and scenario-based testing are designed to support this scale, the cognitive and organizational complexity of maintaining such extensive rulebases cannot be eliminated entirely [37, 38].

Ensuring consistency across interconnected rule sets, preventing logic duplication, and managing change control remain significant engineering challenges. Future development of higher-level abstraction mechanisms, modularization frameworks, and automated rule analysis tools will be essential to support the long-term sustainability of Rule-based decision tools in business environments.

LLM prompting requires governance

Although LLMs can dramatically accelerate rule model creation by translating natural language policies into candidate rule structures, their outputs remain sensitive to prompt phrasing, model version, and context. Without guardrails, an LLM may introduce subtle misinterpretations of policy language or fail to capture critical edge conditions.

As a result, LLM use in rule based system workflows must operate within a formal governance framework that includes versioned prompts, validation checkpoints, reproducibility controls, and mandatory human review. Organisations should treat LLM-generated rules as suggestions rather than executable policy until they have passed automated linting, regression testing, and expert verification. Governance around LLM prompting is therefore essential to prevent unintentional policy drift and ensure regulatory compliance.

Human review remains essential

Despite the deterministic semantics and comprehensive explanation capabilities, the correctness of a rule based model ultimately depends on the human experts who define and validate the underlying policies. Automated evaluation can detect structural issues in rules, but it cannot determine whether the encoded policy is itself fair, lawful, or accurate. Domain experts must review proposed rule changes, validate scenario outcomes, and assess alignment with business intent and regulatory requirements.

Furthermore, because LLM-generated rules are not inherently trustworthy, expert oversight is critical to ensure that natural language ambiguity does not propagate into the decision logic. Rule-Based systems therefore enhance, but does not replace, the human governance processes required for safe and responsible AI deployment in high-stakes domains.

Business integration complexity

Deploying rule based system in production requires integration with existing business systems such as loan origination platforms, underwriting engines, CRM systems, and compliance monitoring workflows. These integrations must accommodate data ingestion pipelines, identity and access management, monitoring infrastructure, and version-controlled deployment environments. Additionally, organizations may need to align decision outputs with downstream audit, case-management, and reporting systems to satisfy internal and external regulatory requirements. Ensuring seamless bidirectional interaction between the rule based system, the LLM services, and operational databases can be non-trivial, particularly in environments with strict security and governance constraints. Successful business adoption therefore depends on careful architectural planning, robust DevOps practices, and comprehensive testing of end-to-end decision flows.

The Rule Authoring Bottleneck

While rule-based systems provides deterministic execution, the cost of rule authoring remains substantial. Domain experts are needed to author and validate rule sets for moderately complex policies. This can represent a significant upfront investment compared to supervised ML approaches, which may achieve comparable accuracy with

labelled data alone. Organisations must weigh this authoring cost against the benefits of interpretability and governance.

LLM-Generated Rules Are Not Automatically Correct

A critical tension exists in the hybrid approach: we argue LLMs cannot be trusted for decision execution due to hallucination, yet we propose using LLMs for rule authoring. This is not a contradiction but rather a risk-reduction strategy. LLM errors in rule authoring are detectable through validation, testing, and expert review before deployment, whereas LLM errors in runtime decision-making are not. However, this does not eliminate the risk of LLM-introduced policy misinterpretations, and organizations must implement rigorous validation of workflows.

Comparison with Modern Interpretable ML

Our evaluation does not include comparisons with other ML based approaches such as gradient boosted decision trees (XGBoost, LightGBM), Generalized Additive Models (GAMs), or modern rule-learning approaches (RuleFit, Skope-rules). These methods offer different tradeoffs between interpretability, performance, and governance requirements. Rule-based systems are most appropriate when: (1) explicit policy encoding is required by regulation, (2) decisions must be reproducible across system versions, or (3) domain experts require direct control over decision logic. For applications where learned patterns are acceptable, interpretable ML may be more cost-effective.

7.3 Broader Impacts

The adoption of a Rule-based system as a symbolic substrate beneath LLM interfaces has the potential to significantly improve the transparency, accountability, and trustworthiness of AI systems used in socially consequential domains. By ensuring that all operational decisions are computed deterministically and accompanied by faithful, structurally grounded explanations, Rule-based systems reduce the risk of opaque model behavior and enables more robust oversight by auditors, regulators, and internal governance teams. This stands in contrast to existing black-box or LLM-only systems, which may obscure sources of error, embed hidden biases, or provide explanations that diverge from actual reasoning processes. A Rule-based system therefore acts not only as a technical reliability layer but also as an institutional governance aid, making complex AI systems more accessible to non-technical stakeholders involved in compliance, policy review, and risk management.

However, having a Rule-based system does not eliminate broader ethical risks inherent to automated decision-making. The accuracy and fairness of the system ultimately depend on the quality of the policies, data, and domain knowledge encoded into such a system as rules, and these may reflect historical inequities or incorrect assumptions.

While a Rule-based system increases the inspectability of such issues, it cannot resolve them on its own; responsible deployment requires rigorous fairness analysis, legal review, and continuous monitoring for disparate impact.

Furthermore, the use of LLMs in rule authoring introduces new challenges, including the potential importation of societal biases or misinterpretation of policy language.

Although a Rule-based systems deterministic execution prevents these issues from directly appearing in operational decisions, organizations must employ appropriate validation workflows to ensure that LLM-generated rules align with domain standards and ethical norms. Overall, the SESL + LLM hybrid approach provides a strong foundation for safer business AI, but it must be embedded within a comprehensive governance strategy to realize its full societal benefits.

7.4 Decision Framework: When to Use a Rule-based System vs. Alternatives

A Rule-based system is not universally optimal. We recommend when:

Strong Indicators :

- Regulatory requirements mandate explicit, auditable decision logic
- Policies change frequently and must be traceable to business decisions
- Domain experts must directly control decision rules
- Decisions must be reproducible across system versions for compliance
- Post-hoc explanations are insufficient (pre-hoc interpretability required)

Indicators Against :

- Optimal patterns are unknown and must be learned from data
- Decision boundaries are highly complex and non-linear
- Authoring cost exceeds model training and validation cost
- Domain expertise is limited or unavailable
- Rapid prototyping and iteration are priorities

Alternative Approaches to Consider:

LLM's and Rule-based Systems will not solve the broad range of challenges alone. When considering how to approach a solution there are more methods and approaches which should be considered :

- Supervised ML with interpretability constraints: When patterns must be learned but some interpretability is needed (use: constrained trees, GAMs, RuleFit)
- Traditional ML + explanation layer: When black-box performance is acceptable and post-hoc explanations suffice (use: XGBoost + SHAP)
- Hybrid symbolic-ML: When some rules are explicit and others learned (use: rule-based preprocessing + ML)

7.5 Ethical Considerations

Potential Positive Impacts:

- Increased transparency in automated decision-making
- Reduced risk of opaque algorithmic harm
- Enhanced regulatory compliance and public trust
- Preservation of institutional knowledge in interpretable form

Potential Negative Impacts:

- False sense of security: Interpretable ≠ fair or correct
- Regulatory arbitrage: Organizations might use SESL for compliance theater while actual decisions remain opaque
- Automation of discrimination: Rules can encode bias as easily as ML models
- Reduced innovation: Explicit rule systems may discourage exploration of novel decision patterns

Deployment Recommendations:

1. SESL should complement, not replace, fairness auditing
2. Organizations must test for disparate impact regardless of interpretability
3. Expert review should include diverse perspectives, not just technical validation
4. Continuous monitoring for unintended consequences is essential

Dual-Use Concerns:

While SESL is designed for legitimate business governance, it could potentially be misused for:

- Automating discriminatory decisions with plausible deniability
- Encoding ethically questionable policies in "interpretable" form
- Circumventing algorithmic accountability requirements through technical compliance

We advocate for SESL deployment only within comprehensive AI governance frameworks that prioritize fairness, accountability, and human oversight.

8. Conclusion

SESL demonstrates that deterministic symbolic execution can address specific limitations of LLM-only approaches in high-stakes business AI particularly hallucination, non-determinism, and ungovernable explanations.

Our evaluation shows that hybrid LLM-SESL architectures achieve perfect determinism and eliminate fabricated justifications, though at the cost of increased authoring effort and reduced flexibility compared to learned models.

Rule-based Systems such as SESL are not universal solutions. It is most appropriate for process and regulated domains where explicit policy encoding, direct expert control, and pre-hoc interpretability are required or strongly preferred. Organisations should evaluate Rule-based approaches against modern interpretable ML alternatives (gradient boosted trees, GAMs, rule-learning systems) based on their specific governance requirements, available expertise, and cost constraints.

Key contributions

1. A modern expert system language with business tooling
2. Demonstration that symbolic substrates can mitigate LLM risks
3. A validation framework for LLM-generated rules
4. Empirical evidence on the determinism-flexibility tradeoff

Future work should focus on:

- Real-world deployment studies with cost-benefit analysis
- Formal verification methods for rule sets
- Automated bias detection and remediation tooling
- Hybrid approaches combining Rule-based Systems with learned components

We believe Rule-based systems such as SESL represents a practical path toward trustworthy AI for specific business contexts, but acknowledge it is one tool among many in the broader interpretable AI landscape.

9. Data and Code Availability

Synthetic Test Data: The 300 synthetic scenarios used in our evaluation are available at www.sesl.ai/paper or upon request to reviewers. These scenarios contain no sensitive information and can be freely reproduced.

SESL Engine: The SESL runtime implementation is made available for research purposes under a non-commercial license at <https://www.sesl.ai>. Academic researchers may request full source code access for independent verification.

Evaluation Scripts: Python scripts for computing evaluation metrics (determinism, fidelity, hallucination detection) are available at www.sesl.ai/paper or in supplementary materials.

LLM Prompts: All prompts used for LLM-only baseline and rule generation are documented in supplementary materials to enable replication.

Reproducibility Note: LLM outputs may vary across API versions even with fixed parameters. Our evaluation used chatGPT5.1. Independent researchers attempting replication should expect minor numerical variation but qualitatively similar results.

References

- [1] Ji, Z., Lee, N., Frieske, R., Yu, T., Su, D., Xu, Y., Ishii, E., Bang, Y. J., Madotto, A., & Fung, P. (2023). Survey of Hallucination in Natural Language Generation. *ACM Computing Surveys*, 55(12), 1-38.
- [2] Huang, L., Yu, W., Ma, W., Zhong, W., Feng, Z., Wang, H., Chen, Q., Peng, W., Feng, X., Qin, B., & Liu, T. (2023). A Survey on Hallucination in Large Language Models: Principles, Taxonomy, Challenges, and Open Questions. *arXiv preprint arXiv:2311.05232*.
- [3] Zhang, Y., Li, Y., Cui, L., Cai, D., Liu, L., Fu, T., Huang, X., Zhao, E., Zhang, Y., Chen, Y., Wang, L., Luu, A. T., Bi, W., Shi, F., & Shi, S. (2023). Siren's Song in the AI Ocean: A Survey on Hallucination in Large Language Models. *arXiv preprint arXiv:2309.01219*.
- [4] Turpin, M., Michael, J., Perez, E., & Bowman, S. R. (2024). Language Models Don't Always Say What They Think: Unfaithful Explanations in Chain-of-Thought Prompting. In *Proceedings of the 37th Conference on Neural Information Processing Systems (NeurIPS 2023)*.
- [5] Lanham, T., Chen, A., Radhakrishnan, A., Steiner, B., Denison, C., Hernandez, D., Li, D., Durmus, E., Hubinger, E., Kernion, J., Lukošiūtė, K., Nguyen, K., Cheng, N., Joseph, N., Schiefer, N., Rausch, O., McCandlish, S., Kundu, S., Kadavath, S., ... Perez, E. (2023). Measuring Faithfulness in Chain-of-Thought Reasoning. *arXiv preprint arXiv:2307.13702*.
- [6] Wiegreffe, S., & Marasović, A. (2021). Teach Me to Explain: A Review of Datasets for Explainable Natural Language Processing. In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*.
- [7] OECD. (2019). *Recommendation of the Council on Artificial Intelligence, OECD/LEGAL/0449*. OECD Legal Instruments.
<https://legalinstruments.oecd.org/en/instruments/OECD-LEGAL-0449>

- [8] National Institute of Standards and Technology (NIST). (2023). *Artificial Intelligence Risk Management Framework (AI RMF 1.0)*. U.S. Department of Commerce.
- [9] Goodman, B., & Flaxman, S. (2017). European Union Regulations on Algorithmic Decision-Making and a "Right to Explanation". *AI Magazine*, 38(3), 50-57.
- [10] Wachter, S., Mittelstadt, B., & Floridi, L. (2017). Why a Right to Explanation of Automated Decision-Making Does Not Exist in the General Data Protection Regulation. *International Data Privacy Law*, 7(2), 76-99.
- [11] Rudin, C. (2019). Stop Explaining Black Box Machine Learning Models for High Stakes Decisions and Use Interpretable Models Instead. *Nature Machine Intelligence*, 1(5), 206-215.
- [12] De Raedt, L., Dumančić, S., Manhaeve, R., & Marra, G. (2020). From Statistical Relational to Neuro-Symbolic Artificial Intelligence. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence (IJCAI-20)* (pp. 4943-4950).
- [13] Besold, T. R., d'Avila Garcez, A., Bader, S., Bowman, H., Domingos, P., Hitzler, P., Kühnberger, K.-U., Lamb, L. C., Lowd, D., Lima, P. M. V., de Penning, L., Pinkas, G., Poon, H., & Zaverucha, G. (2017). Neural-Symbolic Learning and Reasoning: A Survey and Interpretation. *arXiv preprint arXiv:1711.03902*.
- [14] Garcez, A. d., Gori, M., Lamb, L. C., Serafini, L., Spranger, M., & Tran, S. N. (2019). Neural-Symbolic Computing: An Effective Methodology for Principled Integration of Machine Learning and Reasoning. *Journal of Applied Logics*, 6(4), 611-632.
- [15] Chen, Y., Zhong, R., Zha, S., Karypis, G., & He, H. (2022). Meta-Learning via Language Model In-context Tuning. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)* (pp. 719-730).
- [16] Kline, D. M., & Walters, D. J. (2021). Machine Learning in Credit Risk Modeling: Efficiency Comes at a Cost. *The Journal of Risk Model Validation*, 15(1), 91-110.
- [17] Caruana, R., Lou, Y., Gehrke, J., Koch, P., Sturm, M., & Elhadad, N. (2015). Intelligible Models for HealthCare: Predicting Pneumonia Risk and Hospital 30-day Readmission. In *Proceedings of the 21st ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (pp. 1721-1730).
- [18] Selbst, A. D., & Powles, J. (2017). Meaningful Information and the Right to Explanation. *International Data Privacy Law*, 7(4), 233-242.
- [19] Jackson, P. (1998). *Introduction to Expert Systems* (3rd ed.). Addison-Wesley.
- [20] Financial Conduct Authority (FCA). (2023). *Consumer Credit Sourcebook (CONC)*. FCA Handbook. <https://www.handbook.fca.org.uk/handbook/CONC/>

- [21] OECD. (2017). *International VAT/GST Guidelines*. OECD Publishing, Paris.
<https://doi.org/10.1787/9789264271401-en>
- [22] Marcus, G., & Davis, E. (2020). *Rebooting AI: Building Artificial Intelligence We Can Trust*. Pantheon Books.
- [23] Chao, P., Robey, A., Dobriban, E., Hassani, H., Pappas, G. J., & Wong, E. (2023). Jailbreaking Black Box Large Language Models in Twenty Queries. *arXiv preprint arXiv:2310.08419*.
- [24] SESL Engine Source Code. (2025). rule_engine.py. SESL documentation at
<https://www.sesl.ai>
- [25] SESL CLI User Guide. (2025). <https://www.sesl.ai>
- [26] SESL Language/User Guide. (2025). <https://www.sesl.ai>
- [27] SESL Product Summary Guide. (2025). <https://www.sesl.ai>
- [28] ISO. (2018). *ISO 31000:2018 Risk Management, Guidelines*. International Organization for Standardization.
- [29] European Parliament and Council. (2016). Regulation (EU) 2016/679 (General Data Protection Regulation). *Official Journal of the European Union*, L119, 1-88.
- [30] Lipton, Z. C. (2018). The Mythos of Model Interpretability: In Machine Learning, the Concept of Interpretability is Both Important and Slippery. *Queue*, 16(3), 31-57.
- [31] Park et al. (2023) *Generative Agents: Interactive Simulacra of Human Behaviour* (UIST)
- [37] Brachman, R. J., & Levesque, H. J. (2004). *Knowledge Representation and Reasoning*. Morgan Kaufmann.
- [38] Vanthienen, J., & Mues, C. (2006). *Audit and Verification of Business Rule Models*. *Expert Systems with Applications*, 30(4), 570–582.

Acknowledgments and Disclosures

Conflict of Interest: The author is the creator of SESL and is the Founder of Oblongix Ltd, which owns the SESL intellectual property. SESL is released under an open-source non-commercial license.

Funding: This research was conducted as part of product development at Oblongix Ltd. No external funding was received.

Data Availability: Synthetic test scenarios used in evaluation are available at www.sesl.ai/paper/data or upon request for academic review purposes.

Code Availability: The SESL engine implementation is available under license by contacting SESL at <https://www.sesl.ai> for academic use and independent verification.