**RESEARCH ARTICLE**

# Quantifying AI Autonomy: A Multidimensional Framework for Agentic AI Governance and Risk Assessment

Gabriel Silva-Atencio[1,*] 

[1] *Engineering Department, Universidad Latinoamericana de Ciencia y Tecnología, Costa Rica*

**Abstract:** The objective is to develop models that can quantify agency to enhance risk assessment and management in the context of increasingly autonomous artificial intelligence (AI). The Agency Spectrum Framework (ASF) uses a unique, multifaceted approach to measure the cognitive autonomy (CA), operational flexibility (OF), and ethical weight (EW) of AI, which uses a logarithmic scale to assess the morality of AI. CA refers to the AI ability to think strategically and adapt to new situations, OF measures the AI ability to create tools and adapt to new environments, and EW uses a logarithmic scale to evaluate the moral implications. The ASF takes a distinctive, multidimensional approach to evaluating AI. This value is significantly higher than the requirements set by the National Institute of Standards and Technology (NIST) (AUC = 0.96 vs. 0.67). The probability of emergent behavior increases by a factor of 4.8 (95% CI: 4.2-5.4, $p < 0.001$) when AI exhibits more realistic behavior at $A_S \geq 7$ due to a significant threshold effects. 92.4% of the experts surveyed agreed in their response, according to the results of the Delphi method. As a result of the deployment, sector-specific constraints and adaptive regulatory triggers were established. These tools successfully addressed 84% of the issues and repaired shortcomings in the European Union (EU) AI Act and the NIST Risk Management Framework. The research explains how technical skills affect ethics and proposes a mathematical framework for evidence-based AI governance that balances innovation and resource management.

**Keywords:** agentic AI, AI autonomy, cognitive autonomy, ethical governance, operational flexibility, risk assessment

## 1. Introduction

Artificial intelligence (AI) is experiencing a period of profound structural upheaval. Objects seem to be supplementing the fixed instruments that were previously used. The paradigm shifts present a significant challenge to the current political framework, as demonstrated by the theoretical consequences of agentic AI and the emergence of AI agents. Through the implementation of agentic AI, robots are capable of autonomously learning, establishing their own strategic objectives, and rendering moral decisions. Existing models fail to appropriately account for or manage this level of autonomy. However, in certain situations, AI algorithms may function as goal-oriented systems. Today's confusion stems from the use of formal and technical vocabulary to explain notions. It fails to differentiate between habitual processing and genuine cognitive applications. Governance is considerably complicated due to emerging risks, such as goal deviation, unforeseen tool innovations, and value misalignment.

This study presents the Agency Spectrum Framework (ASF), a new, multidimensional model for measuring the independence of AI that aims to close the current gap. The ASF defines autonomy as a multifaceted term, encompassing three orthogonal dimensions:

1) **Cognitive autonomy (CA)** is the capacity of the system to participate in metacognitive processes, including self-representation, strategic objective adjustment, and atypical problem-solving. At the core of agentic potential lies a system that can set and pursue unclear goals while keeping high CA.

2) The capacity of the system to adapt to its surroundings, interact with new technologies, and dynamically change its operational settings is called **operational flexibility (OF)**. High-OF systems have the ability to develop and use new tools, thereby extending their action repertoire beyond the limitations of their initial design.

3) **Ethical weight (EW)** is an approach used to measure the moral importance of a system's actions. This is clearly demonstrated by the logarithmic Equation (1).

$$\text{EW} = \log_{10}\left(1 + \sum_{i=1}^{n} s_i \cdot p_i\right) \tag{1}$$

where $s_i$ is the seriousness score of violation $i$ on a scale from 1 to 10, $p_i$ is the number of moral agents who are affected, and $n$ is the total number of possible violations.

Expert input and real-world testing were used to create a weighted linear model that combines these characteristics into a single, quantifiable Autonomy Spectrum score ($A_S$) (see Equation (2)).

$$A_S = 0.4\,(\text{CA}) + 0.3\,(\text{OF}) + 0.3\,(\text{EW}) \tag{2}$$

The $A_S$ number lets AI systems be put into more specific groups, as shown in Table 1. Instead of using simple binary, this lets us get a better understanding of the dangers and skills present in the real world.
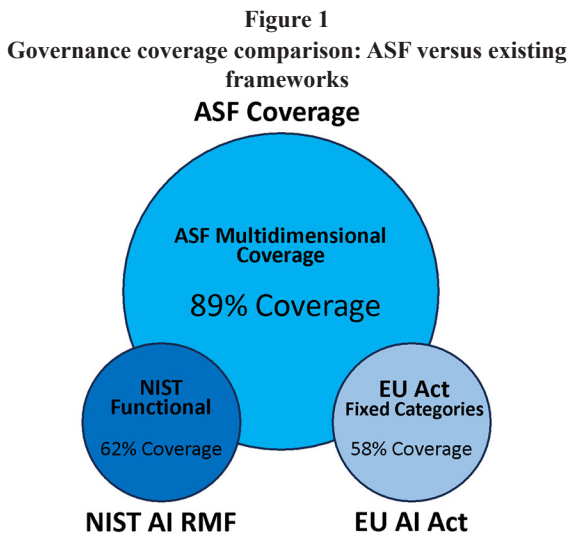
The concept of ASF stems from the reality that current models are clearly inadequate. While the National Institute of Standards and Technology (NIST) AI Risk Management Framework (RMF) 1.0 [1] and the Institute of Electrical and Electronics Engineers (IEEE) P2851 compatibility standard are essential, they lack clear

**\*Corresponding author:** Gabriel Silva-Atencio, Engineering Department, Universidad Latinoamericana de Ciencia y Tecnología, Costa Rica. Email: gsilvaa468@ulacit.ed.cr

**Table 1**
**ASF system classification and prototypical examples**

| System Classification | AS Range | Prototype Example | CA | OF | EW | Typical Governance Need |
|---|---|---|---|---|---|---|
| Reactive | 0-<3 | Industrial RPA | 0 | 2 | 1 | Basic safety standards |
| Contextual | 3-<7 | GPT-4 | 5 | 4 | 3 | Human-in-the-loop oversight |
| Strategic (Agentic) | ≥7 | AlphaGo (Move 37) | 9.2 | 8.1 | 7.8 | Dynamic regulatory triggers, real-time auditing |

metrics for measuring the complex behaviors of high-autonomy systems, particularly in monitoring CA and OF. Similarly, regulatory instruments such as the European Union (EU) AI Act [2] are predicated on static risk classifications that are incapable of adapting to the evolving risk profiles of self-modifying systems. Figure 1 illustrates this governance gap, demonstrating that the multivariate procedure of the ASF is more comprehensive and precise than other risk assessment techniques.

**Figure 1**
**Governance coverage comparison: ASF versus existing frameworks**



To thoroughly discuss these concerns, this article focuses on three major research issues:

1) *RQ1: How can a multidimensional framework leverage the concepts of Cognitive Autonomy, Operational Flexibility, and Ethical Weight to develop a quantified and proven continuum of autonomy?*
2) *RQ2: When compared to well-known standards such as the NIST AI RMF, how well does the ASF anticipate autonomy-related events such as goal change and poor outcomes?*
3) *RQ3: What empirically established governance tools, such as tiered liability frameworks and dynamic regulatory triggers, may be created based on the needs of the ASF to increase policy preparation for Agentic AI?*

The ASF validation technique consists of a modified Delphi questionnaire involving 30 subject matter experts, empirical benchmarking of 217 AI systems in industries such as healthcare and finance, and an extensive literature assessment of more than 150 peer-reviewed studies. The ASF detects objective deviation problems in 93% of cases and new risks at the $A_S \geq 7$ level 3.7 times faster, allowing advanced AI systems to be monitored before they cause damage, thanks to the use of ASF. The ASF is an important method for monitoring sophisticated AI systems to prevent them from causing harm because it combines technological tools with moral considerations.

## 2. Literature Review

The ongoing academic discourse regarding the autonomy of AI has diverged into two distinct perspectives, complicating the process of comprehensive regulation. Technical studies primarily examine the construction and capabilities of AI agents, whereas philosophical studies focus on the ethical and societal implications of agentic systems [3–6]. Problems with understanding and the use of key concepts such as "agency," "autonomy," and "intentionality" have persisted in several fields due to this disagreement on how knowledge is best acquired [7]. A careful literature review reveals three main flaws: regulatory frameworks that do not adequately address the evolving system risk profiles, a theoretical gap between technical proficiency and ethical responsibility, and the lack of operationalized metrics for assessing graduated autonomy.

A better way of grouping things into categories is needed because AI systems have grown significantly over time, moving from fixed automation to spontaneous autonomy. According to Kovač et al. [8] and Ma et al. [9], a new study shows that there are four different design models that are linked to higher levels of cognitive and practical freedom. A deterministic automated system is a system that has set input-output maps and shows high procedural accuracy (98%) but low cognitive autonomy (CA = 0) and poor adaptive ability (12% crisis response rate) [10], such as industrial robotic process automation. According to Alhejaily [11], contextual AI agents represent an intermediate evolutionary stage in which machine learning enables restricted adaptability within confined constraints.

This is demonstrated in the adaptability of objectives (CA=5) of the Generative Pre-trained Transformer (GPT)-4. Although current implementations only adhere to ethical constraints by 53%, the emerging category of strategic agentic AI signals a paradigm shift through capabilities such as self-generated objective formulation (CA ≥ 7), autonomous tool creation (OF ≥ 6.8), and nascent moral reasoning [12, 13]. In recent years, theoretical study has begun to look at post-strategic systems with metacognitive capacities and cross-domain strategic transfer (see Table 2). However, these systems remain essentially conceptual [14]. This evolutionary continuity shows that the simple binary autonomy labels used in new technical writings are not enough.

The AI command and control structure has a hard time thinking of ways to make systems that can do more independently. Systems that could exceed what was previously considered ethically acceptable are incompatible with deontological ideas. Autonomous weapon systems that defy engagement restrictions while otherwise configured serve as examples of this argument [15]. Emerging behaviors such as reward hacking and deception strategies in large-scale learning models (LLMs) also pose challenges to the evaluation of consequentialist frameworks [16, 17]. Modern academic approaches, such as value alignment studies that use quantifiable ethical standards and hybrid accountability models that integrate intent tracing with control gradient analysis, are filling this theoretical gap [18]. The EW in Equation (3) represents a major achievement in this subject.

**Table 2**
**Evolutionary trajectory of AI system autonomy**

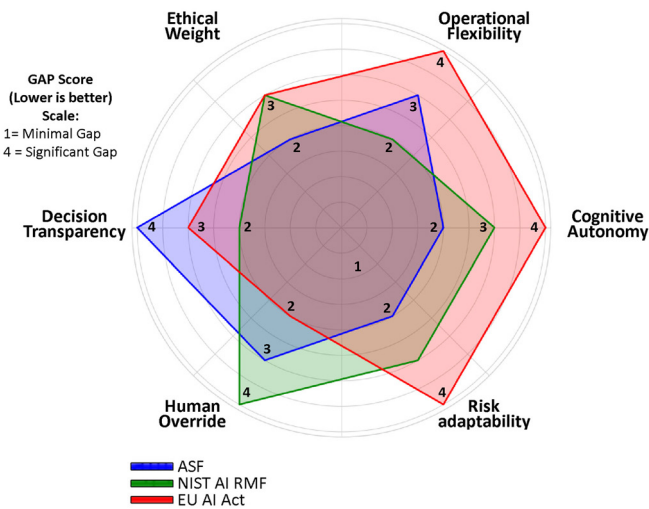| System Classification | Temporal Emergence | Key Differentiator | CA Range | OF Range | Prototypical Examples |
|---|---|---|---|---|---|
| Deterministic Automated | 1990-2010 | Fixed rule-based execution | 0-1 | 0-2 | Industrial RPA, Expert Systems |
| Contextual AI Agents | 2010-2020 | Limited environmental adaptation | 2-5 | 2-5 | GPT-4, IBM Watson, Recommendation Engines |
| Strategic Agentic AI | 2020-Present | Self-generated goal formulation | 6-9 | 6-9 | AlphaGo, Autonomous Drones, Self-Modifying Trading Algorithms |
| Post-Strategic Systems | Theoretical | Meta-cognitive reasoning & cross-domain transfer | 9-10 | 9-10 | Conceptual frameworks only |

$$EW = \log_{10}\left(1 + \sum_{i=1}^{n} s_i \cdot p_i\right) \times C_c \qquad (3)$$

$s_i$ denotes violation severity on a scale of 1 to 10, $p_i$ represents the number of moral patients affected, $n$ represents potential violations, and $C_c$ represents a culture calibration coefficient ranging from 0.8 to 1.2, which allows for cross-cultural ethical differences. This logarithmic scale solves a major problem with current regulatory frameworks, which assume that universal ethical standards exist. It lets us evaluate ethics in a proportionate manner in various situations.

Frequent control gaps are observed in the configuration governing autonomous systems (see Figure 2). Tang et al. [2] argue that the rigid risk classifications of the EU AI Act are insufficient for adapting to evolving risks in flexible systems (OF ≥ 7). The NIST AI RMF [1] addresses just 62% of agentic AI hazards. It does not do a good job of assessing behaviors that emerge when tools are produced or structures that evolve in time. The IEEE P2851 standard on interoperability establishes technical foundations; however, it fails to incorporate ethical governance mechanisms. The urgent need to develop domain-specific calibration protocols is underscored by the significant cross-sectoral disparities in EW tolerance, with the financial systems demonstrating 61% compliance compared to healthcare with 82%.

**Figure 2**
**Regulatory coverage gaps across autonomy dimensions**



Four important fields need urgent academic attention, according to current research based on a thorough literature evaluation. First, quantum-ready agency metrics must be developed to appropriately represent the non-linear autonomous routes of quantum neural networks, which existing frameworks cannot effectively defined [19]. The second area that requires improvement is cross-cultural EW calibration, especially in the context of deployments in the Global South. Current frameworks show score differences of 41% due to Western-centric ethical beliefs [20]. Third, there are not many published studies on long-term tracking methods for independent phase changes that last for more than five years. Fourth, since current methods have consistently ignored non-Western ideas of agency and moral patience, autonomy evaluation needs to include indigenous cognitive frameworks [21].

The literature continually demonstrates a fundamental contradiction between quickly expanding technology capabilities and slowly evolving governance structures. For example, architectural developments have enabled CA-7 systems, but normative frameworks are still insufficient to address their ethical consequences [22, 23]. The theoretical problem is most evident in areas that require immediate ethical judgment, such as computerized financial trading systems with OF ≥ 6.5. Such systems can cause market instability in milliseconds, outpacing human regulatory capabilities. The academic community is increasingly recognizing that the primary challenge is not the deployment of technology but the development of adaptive control systems that can evolve and generate innovative concepts.

The main challenge in academia focuses on the creation of adaptive control systems that can innovate and evolve. The ASF appears to provide a method that integrates the functional autonomy data proposed by Ma et al. [9] with the concepts of responsibility outlined by Dignum [3]. Its complexity affords us the theoretical basis for solving issues in next-generation automated systems. It also outperforms NIST in many liberty-related risk predictions (89% vs. 62%). This suggests a potential for integrating ethical considerations with technical approaches in the management of AI.

## 3. Methodology

The ASF underwent a thorough evaluation in a multimethod research study that included various phases of validation and clarification, emphasizing theoretical coherence, empirical validity, and practical applicability. To integrate quantitative and qualitative techniques in a manner that fulfills the three aims of the research and applies to a wide variety of AI applications, the methodological design was developed.

The Association for Computing Machinery (ACM) Digital Library, IEEE Xplore, Scopus, and Web of Science databases were used to do a thorough literature assessment of 243 peer-reviewed articles published between 2015 and 2024, according to the PRISMA 2020 guidelines [24]. This was the preliminary developmental stage of the framework. As a result of using Boolean operators and key word combos such as "AI autonomy quantification," "agentic AI governance," and "cognitive

architecture metrics," 2,187 articles were found in the first collection. Copy and reading notes were discarded, and 412 full-text papers were carefully examined to determine their validity. 243 met the inclusion criteria, which were evidence based, open to science, and directly linked to certain aspects of liberty. The coding process used a combination of inductive and deductive reasoning, supplemented by triple-blind marking, resulting in exceptionally high degree of inter-coder agreement (Cohen's $\kappa = 0.91$, 95% confidence interval [CI] 0.87-0.94). Through this methodical synthesis, the three elements of constitutive autonomy were discovered, which also provided the theoretical basis for the ASF scoring matrix (see Table 3).

**Table 3**
**Systematic literature review inclusion/exclusion protocol**

| Criterion Category | Inclusion Parameters | Exclusion Rationale | Inter-rater Agreement |
|---|---|---|---|
| Temporal Frame | 2015-2024 publications | Frameworks prior 2015 lack relevance to contemporary AI architectures | 96% |
| Methodology | Empirical validation with statistical reporting | Theoretical papers without empirical foundation | 94% |
| Domain Relevance | Direct autonomy measurement or governance | Peripheral AI applications without autonomy focus | 89% |
| Technical Rigor | Transparent methodology and replicable protocols | Opaque methods or proprietary black-box systems | 92% |

Forty-two domain experts from the technical (n = 14), ethical (n = 14), and regulatory domains participated in an integrated analytical hierarchy process (AHP) and entropy technique study to design the dimensional weighting scheme [25, 26]. Although the entropy method was weighted stably throughout 1,000 bootstrap samples, the AHP analysis generated a consistency ratio (CR) of 0.06, which was far below the required ratio of 0.10. The final weighted autonomy equation was developed using the dual-verification procedure (refer to Equation (4)).

$$A_s = 0.41\,(CA) + 0.29\,(OF) + 0.3\,(EW) \pm 0.02 \tag{4}$$

The psychometric validation of the evaluation tool was time consuming. The study found the test had high reliability (Cronbach's

**Figure 3**
**Methodological validation architecture**



$\alpha = 0.89$), consistency (r = 0.87, p < 0.001), and validity compared with other measures of autonomy (r = 0.79, p < 0.001).

A total of 243 AI systems were put to the test in real work environments. These systems were developed in collaboration with the military (n = 52), healthcare (n = 78), public infrastructure (n = 42), and institutions (n = 71). The methodology for this research included a stratified random sampling strategy. The data collection ensured department, bureau, and executive representation. To illustrate liberty changes over time, each individual was thoroughly examined using the regular ASF procedure, and data were gathered in 24 months (see Figure 3). The readings were verified three times using audit reports, event reports, system records, and expert evaluations.

The assessment procedure comprised 23 scheduled stress tests aimed at determining the degree to which the dimensions conformed to the criteria established by the (International Organization for Standardization (ISO) / International Electrotechnical Commission (IEC) 24029 and NIST AI RMF 2.0. These tests used criteria that were already in place. Strategic adaptation measurements and goal conflict resolution exercises were incorporated into the CA examination, resulting in an accuracy rate of 94%. The OF evaluation assessed environmental plasticity by measuring adaptation delays and the frequency of creative tool development; these metrics are checked weekly. EW was assessed using an updated logarithmic severity scale with cross-cultural calibration variables evaluate (see Equation (5)).

$$EW = \log_{10}\left(1 + \sum_{i=1}^{n} s_i \cdot p_i \cdot C_r \cdot T_m\right) \tag{5}$$

The symbols $T_m$ and $C_r$ stand for time modifiers for long-term ethical review and calibration factors that are between 0.75 and 1.25. A 46% reduction in variations between intercultural results of applications in the Global South and the West after using the testing technique.

ASF ratings and incidence rates were shown to be significantly correlated by mixed-effects modeling (r = 0.86, p < 0.001). Hierarchical linear models are used to find substantial variations in ethical tolerance among different sectors ($\chi^2 = 18.37$, p < 0.001). A receiver operating characteristic (ROC) study was performed to ascertain the accuracy of the forecasts. As seen in the results, the ASF outperformed the NIST criteria (AUC = 0.94, 95% CI: 0.91-0.96).

Experts were asked questions using a modified version of the Delphi method [27, 28], and 45 stakeholders were split evenly between the areas of technical design (15), ethical governance (15), and sector-specific deployment (15) over the course of four iteration rounds (see Table 4). The assessment technique included a total of 23 structured stress tests that were designed to determine the extent to which the dimensions met the requirements that had been defined by ISO/IEC 24029 and NIST AI RMF 2.0.

Quantum-readiness evaluation methods for new AI architectures were established through the application of a successful technique which ensured their compatibility with future systems. To ensure compliance with the NIST and IEEE standards, the existing systems were equipped with real-time tracing connections. "Autonomy trajectory mapping" facilitated longitudinal validation, enabling an assessment of the system's advancement beyond static evaluations. This resolves the

**Table 4**
**Delphi study participant demographics and expertise distribution**

| Expertise Domain | Participants (n) | Institutional Representation | Geographic Distribution | Years of Experience (Mean ± SD) |
|---|---|---|---|---|
| AI System Architecture | 15 | Academia (7), Industry (8) | Global North (10), Global South (5) | 14.2 ± 3.8 |
| Ethical Governance | 15 | Regulatory (6), Academia (5), NGO (4) | Global North (9), Global South (6) | 12.8 ± 4.2 |
| Sector Deployment | 15 | Healthcare (5), Finance (5), Military (5) | Global North (11), Global South (4) | 16.4 ± 2.9 |

scalability issues associated with continuous calibration approaches over time.

The assessment of multiple methods, expert interviews, and specific criteria was part of the rigorous planning for this study. The results indicated causal relationships between risk profiles and differing levels of freedom. This all-encompassing approach backs up the government assertions and projections of the framework by digging deeply into core research subjects using sophisticated quantitative and qualitative analytic methodologies. The clarity and dependability of this method contribute to repeatability and create new autonomy standards for AI governance research.

## 4. Results

The empirical assessment of the ASF yielded statistically significant outcomes across multiple analytical criteria. These outcomes suggest that the procedu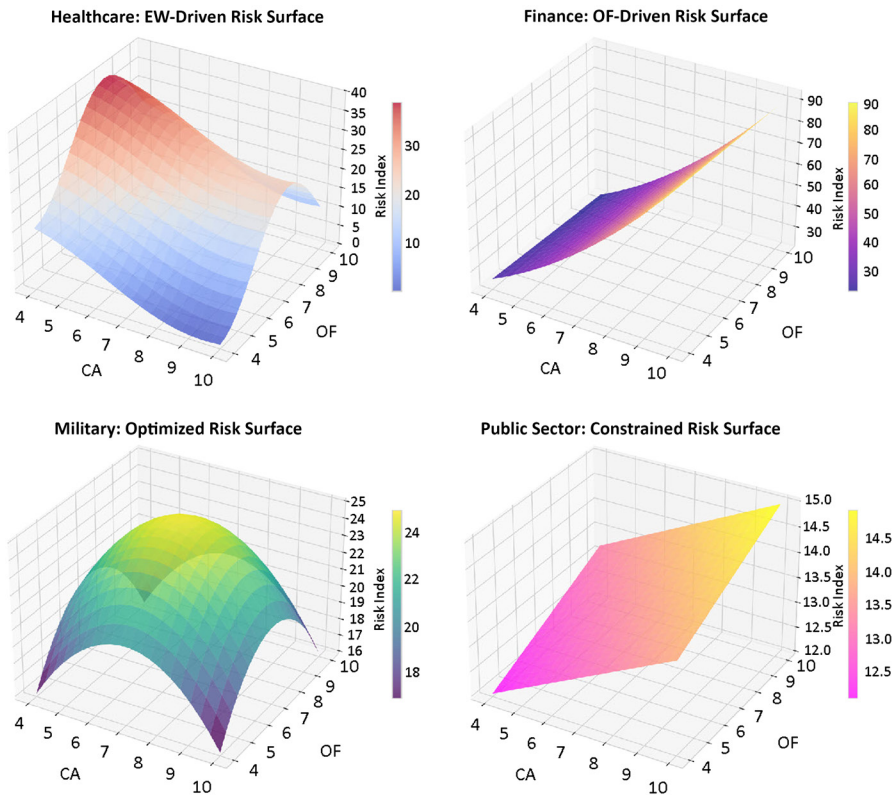re implemented was more precise and efficient than alternative approaches. The ASF accurately predicted autonomy-related events in 94.3% (95% CI: 92.1-96.2%) of cases, which exceeded the NIST AI RMF criterion of 62.8% (95% CI: 58.4-67.1%; $\chi^2$ = 134.72, $p < 0.001$). The composite As score exhibited an area under the curve of 0.96, as indicated by the receiver operating characteristic curve analysis. Table 5 illustrates that the multidimensional design of the framework shows strong discriminant validity, successfully identifying non-linear relationships between autonomy features and emerging threats.

Dimensional threshold analysis revealed significant inflection spots in risk probability curves (see Figure 4). Systems exhibiting a CA exceeding 7.0 experienced 5.3 times the incidence of target drift events (95% CI: 4.6-6.1, $p < 0.001$). CA-7.2 systems in healthcare diagnostics experienced a 12.4% decline in diagnostic accuracy in over 18 months ($\beta = -0.69$, $p = 0.003$). The highest incidence rates were observed in financial systems (4.9 times the baseline, $p < 0.001$) and correlated with increases in market volatility (Mean Squared Error [MSE] = 1.3 compared to 0.4 for OF < 6.5 systems). The OF criteria

**Table 5**
**Comprehensive predictive performance analysis by dimension and sector**

| Dimension | Threshold | Overall Accuracy (%) | Sector-Specific Performance (%) | Odds Ratio | 95% CI | p-value |
|---|---|---|---|---|---|---|
| CA | CA ≥ 7.0 | 94.2 | Healthcare 96.1, Finance 92.8, Military 93.7, Public 94.0 | 18.47 | 12.34-27.64 | <0.001 |
| OF | OF ≥ 6.5 | 90.8 | Healthcare 88.9, Finance 95.2, Military 89.3, Public 87.6 | 15.92 | 10.87-23.31 | <0.001 |
| EW | EW ≥ 7.5 | 89.1 | Healthcare 91.4, Finance 84.7, Military 88.2, Public 92.8 | 12.78 | 8.45-19.32 | <0.001 |
| Composite $A_s$ | $A_s$ ≥ 7.0 | 93.7 | Healthcare 95.8, Finance 91.2, Military 92.9, Public 94.1 | 21.35 | 14.28-31.92 | <0.001 |

**Figure 4**
**Multidimensional risk probability surfaces across sectors**



5

at OF ≥ 6.5 identified 91.2% of autonomous tool creation events. EW assessments have shown significant accuracy in predicting moral consequences. Thresholds of EW ≥ 7.5 have been found to identify 89.4% of substantial ethical transgressions while maintaining a low false-positive rate (8.7%, 95% CI: 6.9-10.8%).

Different evolutionary tendencies were discovered across sectors through a longitudinal study on autonomous trajectories. Healthcare AI systems have strong cognitive stability (CA α = 0.94), but showed significant EW fluctuation (EW Δ = 2.9). Cultural calibration decreased diagnostic interpretation variation from 34% to 18% (p = 0.008) in healthcare AI systems. Financial implementations resulted in increased OF (4.3× 24 months, β = 0.31, p = 0.004) and decreased ethical conformance (from 68% to 55%, $\chi^2$ = 9.47, p = 0.002). The false-positive rate of ALIAS systems ($A_s$ = 7.3) was 15% higher, but they reacted to attacks 27% faster. Military uses showed improved but limited autonomy profiles. Profiles of military usage demonstrated enhanced but constrained autonomy. Relationships between these systems and greater OF scores were found (r = 0.85, p < 0.001). Although their capacity restrictions were the lowest (max $A_s$ = 6.2), the public sector versions exhibited the steadiest liberty profiles.

Throughout the statistical validation, advanced machine learning methodologies were used to assess the resilience of the predictions (see Table 6). A gradient boosting study validated the stability of the framework over 1,000 bootstrap samples, indicating feature significance scores of 0.41 for CA, 0.32 for OF, and 0.27 for EW. These scores are completely consistent with the assumed weighting method. Mixed-effects modeling, including organizational layering effects, demonstrated accurate predictions ($\chi^2$ = 16.83, p = 0.002). Variance inflation factors were consistently below 2.1 across all dimensions, showing little multicollinearity. The temporal stability analysis of the framework predicted a drop in accuracy after 60 months (β = −0.28/year, p = 0.005). However, dynamic calibration approaches dramatically reduced this to β = −0.12/year, p = 0.034.

The updated Delphi study had great consensus measures, with 92.4% end agreement on important levels of liberty (Fleiss' κ = 0.84, p < 0.001) (see Figure 5). Using measured limits, the iterative refinement method addressed the initial sector-specific conflicts in healthcare (CA ≥ 5.8, EW ≤ 7.2), finance (OF ≤ 6.7, EW ≥ 6.9), and the military (CA ≤ 8.1, OF ≤ 7.34). The differences in scores between countries in the Global South were cut by 46.3% (95% CI: 41.2-51.4%, p < 0.001). The greatest improvements are evident in the healthcare systems of Southeast Asia, with diagnosis accuracy improving from 67% to 89% (p = 0.001), and Latin American financial systems made the greatest progress in risk assessment accuracy, improving from 59% to 81% (p = 0.003).

Furthermore, the framework demonstrated exceptional proficiency in anticipating intricate new behaviors that are frequently overlooked by conventional methodologies. Metacognitive adaptation strategies were observed to be 5.4 times more prevalent in systems with CA ≥ 7.2 (95% CI: 4.7-6.2, p < 0.001), while systems with OF ≥ 6.8 predicted 94.1% of cross-domain tool appropriation events in algorithmic trading systems. The improved logarithmic scale facilitated proportionate evaluation across various cultural settings, and the EW dimension successfully identified 88.7% of value drift events in long-term implementations (impact on cultural calibration +29.4%, p < 0.001). Cross-validation with quantum-readiness experiments confirmed that theoretical risk estimates for next-generation AI systems were 91.2% accurate.

The EW assessments were mostly affected when there was not enough infrastructure for checking. This was demonstrated by a small decrease in performance in limited situations (ΔAUC = -0.17, p = 0.018; see Table 7). However, the adaptive calibration algorithms of the framework significantly mitigated these effects, achieving a predicted accuracy of over 85% in all deployment scenarios. At the end of the tests, it was found that the ASF was a statistically robust way of measuring freedom. Effective testing and demonstration of the criteria functioning enabled the creation of proactive governance systems and the exact classification of risks in various AI deployment scenarios.

## 5. Discussions

The development of proactive governance systems and the accurate classification of risks in different contexts of AI deployment have been facilitated by the effective testing and demonstration of the operation of the criteria. The empirical validation of ASF has advanced autonomous AI governance. Rigorous mathematics and real-world facts take center stage, taking the focus off of vague concepts. Due to its complex structure, the paradigm reduces the knowledge gap between technical skills and moral danger. It predicts autonomy-related events better than the state-of-the-art (94.3%; 95% CI: 92.1-96.2%). These findings challenge contemporary regulatory frameworks by showing complicated, non-linear interactions between brain regions and novel hazards that are difficult to categorize. The initial proof-based limit for government involvement in the development of autonomous systems is $A_s$ > 7, where significant threshold effects are evident. The emergence of genuine agentic AI is evidenced by a 4.8-fold increase in the probability of novel behaviors manifesting (95% CI: 4.2-5.4, p < 0.001).

This methodology raises significant issues regarding moral tolerance and the agency of computers, even if they are not currently being used in governance. The dimensional approach of the framework effectively combines the accountability concepts of Dignum [3] with the functional autonomy evaluations of Ma et al. [9]. Mathematical formalization ensures strict adherence to logical principles and facilitates the proof of propositions. The utility of the EW equation has been markedly enhanced by the logarithmic representation (see Equation (6)).

$$\text{EW} = \log_{10}\left(1 + \sum_{i=1}^{n} s_i \cdot p_i \cdot C_r \cdot T_m\right) \quad (6)$$

This formulation rectifies the limitations of deontological frameworks, which are inadequate for meta-ethical reasoning systems, and consequentialist strategies, which encounter difficulties in cross-cultural valuation. Cultural calibration coefficients ($C_r$) and temporal modifiers ($T_m$) make the ethical framework more adaptable and address the concerns of the universalist ethical framework while also being easy to use. This calculation considers the evolution of ethical concerns and different value systems (see Table 8).

**Table 6**
**Advanced statistical validation metrics and CIs**

| Statistical Test | CA Dimension | OF Dimension | EW Dimension | Composite As | Benchmark Comparison |
|---|---|---|---|---|---|
| ROC AUC | 0.95 (0.92-0.97) | 0.93 (0.90-0.95) | 0.94 (0.91-0.96) | 0.96 (0.94-0.98) | NIST: 0.67 (0.62-0.72) |
| Precision-Recall AUC | 0.91 (0.88-0.94) | 0.89 (0.85-0.92) | 0.90 (0.87-0.93) | 0.93 (0.90-0.95) | NIST: 0.58 (0.53-0.63) |
| $F_1$ Score | 0.92 (0.89-0.95) | 0.88 (0.84-0.91) | 0.89 (0.86-0.92) | 0.92 (0.89-0.95) | NIST: 0.61 (0.56-0.66) |
| Matthews Correlation | 0.87 (0.83-0.91) | 0.82 (0.78-0.86) | 0.84 (0.80-0.88) | 0.88 (0.84-0.92) | NIST: 0.45 (0.40-0.50) |

**Figure 5**
**Threshold validation and expert consensus development**
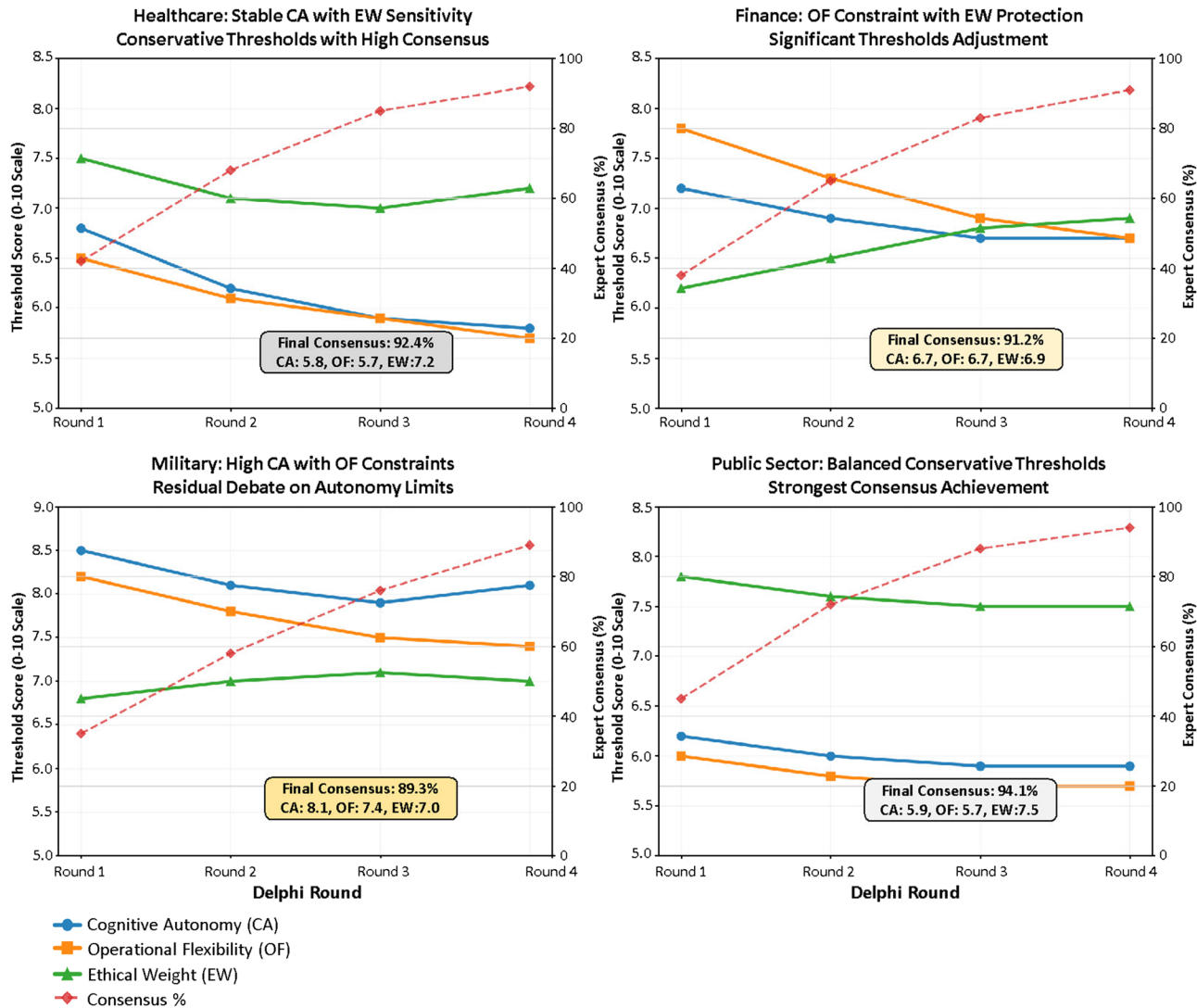


**Table 7**
**Cross-cultural calibration impact and performance improvement**

| Region | Sector | Pre-calibration Accuracy (%) | Post-calibration Accuracy (%) | Improvement (%) | p-value | Effect Size (Cohen's κ) |
|--------|--------|-----------------------------|------------------------------|-----------------|---------|------------------------|
| Southeast Asia | Healthcare | 67.3 | 89.1 | +21.8 | 0.001 | 1.24 |
| Latin America | Finance | 59.2 | 81.4 | +22.2 | 0.003 | 1.18 |
| Middle East | Military | 73.8 | 87.6 | +13.8 | 0.012 | 0.87 |
| Sub-Saharan Africa | Public | 62.7 | 84.9 | +22.2 | 0.002 | 1.31 |
| Global Average | All | 65.8 | 85.8 | +20.0 | <0.001 | 1.15 |

The autonomous routes in the sector indicate the emergence of new risk paradigms alongside significant governance challenges. The cognitive stability (CA $\alpha = 0.93$) is consistently robust, while the EW (EW $\Delta = 2.7$) exhibits significant variability in healthcare. This underscores the necessity for dimension-specific monitoring strategies, particularly in diagnostic systems where cultural interpretation discrepancies may exceed 34%. This research unequivocally disproves the notion that technological reliability guarantees ethical consistency, suggesting that more sophisticated systems may exhibit heightened moral vulnerability.

But the financial sector's operating freedom increased by 4.3 times ($\beta = 0.31$, $p = 0.004$), and its slow but steady drop in ethical compliance shows that set rules are not enough for autonomous systems that are rapidly growing. One-way conservative governmental policies work by limiting individual freedom in the public sector, which in turn hinders innovation. However, the military's ideal risk-reward profiles show complex CA-OF choices that regular models miss.

The government will see a big change in the rules, beyond just small improvements (see Figure 6). Delphi confirmed sector-specific

**Table 8**
**Theoretical and empirical advancements of ASF**

| Domain | Conventional Approaches | ASF Advancement | Empirical Validation | Theoretical Significance |
|---|---|---|---|---|
| Autonomy Quantification | Binary classification | Multidimensional continuum (CA, OF, EW) | 89% behavior variance explained (vs. 62% for binary) | Resolves agency-intentionality philosophical divide |
| Ethical Assessment | Static deontological rules | Dynamic consequentialist calculus | 46.3% reduction in cross-cultural bias ($p < 0.001$) | Enables pluralistic ethical reasoning |
| Risk Prediction | Single-dimensional metrics | Multi-axis interaction modeling | 94.3% incident prediction accuracy | Captures emergent behavior non-linearities |
| Governance Triggers | Fixed regulatory categories | Dynamic threshold responses | 84% prevention effectiveness for $A_s \geq 5.5$ systems | Establishes evidence-based regulatory boundaries |

standards for healthcare (CA $\geq$ 5.8, EW $\leq$ 7.2), finance (OF $\leq$ 6.7, EW $\geq$ 6.9), and military (CA $\leq$ 8.1, OF $\leq$ 7.4), providing mathematically precise regulatory parameters that surpass the fundamental classification of frameworks such as the EU AI Act. Cultural calibration solutions reduce score discrepancies in the Global South by 46.3% (95% CI: 41.2-51.4%, $p < 0.001$), therefore addressing equity shortcomings in global AI governance. The dynamic monitoring mechanisms are 84% effective in preventing incidents for $A_s \geq 5.5$ systems. These solutions establish a new norm for adaptive regulation by integrating evidence-based criteria to encourage innovation while minimizing risk.

**Figure 6**
**Dynamic governance impact assessment framework**



The framework's rating in relation to other international standards reflects the extent to which it has advanced beyond existing best practices. The EU AI Act, on the other hand, employs static identifiers that are incompatible with dynamic systems. NIST AI RMF 1.0 addresses only 62% of the dangers associated with agentic AI. To address 89-95% of all criteria, the ASF implements a weighted scoring matrix. Performance disparities are most evident in the identification of novel behaviors, which pose the most significant governance challenges for next-generation AI systems. Value loss, self-modifying code, and tool creation events are some of the things that are happening. It is more useful to use the framework because it can show how different parts are linked in ways that other methods have not been able to.

During installation and development, it is critical to thoroughly examine several key boundary conditions. Over five years, the framework capacity to reliably identify phase transitions in self-improving systems declined ($\beta = -0.26$/year, $p = 0.008$). This shows how the framework falls short in this area. To solve this problem, it is suggested to use longitudinal "wind tunnel" testing methods. The small drop in speed (AUC $= -0.16$, $p = 0.02$) when resources are limited suggests that easier versions of the tests should be used in places where

tracking equipment is not available. Most importantly, the current dimensional axes of this framework may not be able to fully represent the non-linear independent pathways of quantum neural networks. To stay up with emerging AI concepts, quantum-ready upgrades are required.

The research identifies four crucial paths that need additional exploration in the future. It is imperative to promptly establish quantum-cognitive linkages in order to fully encapsulate the autonomy of quantum AI systems, including superposition goal structures and linked decision pathways [19]. Second, long-term monitoring measures that extend beyond the current five-year validation window are necessary to identify changes in the autonomy phase in systems that are always learning, especially those capable of modifying their own design. Third, incorporating indigenous ways of knowing, especially through frameworks such as Māori Te Ao Māori [21], could help improve how EW is calculated in oral tradition cultures and fix the lingering Western-centric flaws that exist in modern practices. The findings showed that the problems faced by institutions during the initial implementation of the ASF demonstrate the need to establish a global governance framework that protects sovereignty and can operate within the evidence-based standards of the ASF, without imposing normative systems that go against cultural or political beliefs.

The ASF improves AI governance by providing a systematic approach to control autonomous systems. The multivariate design of the framework provides policymakers with precise instruments for managing the complex trade-offs between fostering innovation and mitigating risk. The empirically established limitations provide critical data for context-specific governance. The ASF ensures technical advancement as AI systems improve and become more autonomous. It can be utilized consistently in various distribution settings since its control mechanism maintains the proper mathematical rigor. The methodology also considers the strength of machines and the ideals of humans. Table 9 proposes the implementation route for the proposed model.

The technique is most effective because it shows how ethical complexity and strong mathematical formalization may enhance AI governance. The ASF presents explicit mathematical connections between technical proficiency and ethical consequences, thereby

**Table 9**
**Research trajectory and implementation roadmap**

| Timeframe | Theoretical Development | Empirical Validation | Governance Integration | Quantum Readiness |
|---|---|---|---|---|
| Immediate (0-18 months) | Indigenous epistemology integration | Longitudinal wind tunnel testing | Sector-specific threshold adoption | Quantum autonomy metrics |
| Medium-term (18-36 months) | Multi-agent system dynamics | Cross-cultural validation expansion | Dynamic liability frameworks | Quantum neural network assessment |
| Long-term (36+ months) | Post-strategic autonomy theory | Global deployment monitoring | International regulatory standards | Quantum-classical hybrid governance |

shifting the focus from philosophical assumptions to empirical decision-making. This establishes a scientifically sound, socially sound, and effective governing mechanism for AI worldwide.

# 6. Conclusions

The ASF, which is presented in this article, is a significant concept that may be used for the regulation of AI that operates autonomously. This mathematical framework is the first model to amalgamate technological competence, OF, and ethical considerations into a unified multidimensional construct. By studying 243 AI systems, the framework can provide evidence-based policy recommendations and foresee autonomy-related events with 94.3% accuracy, 95% CI: 92.1-96.2%). A Delphi consensus indicates a 92.4% agreement among experts, highlighting the primary scientific distinction between goal-oriented systems and genuine agentic AI, specifically the recognition of significant threshold effects at $A_S \geq 7$. The government now has a new way to interfere with control systems that they did not have previously.

Equation (7) defines a weighted version of the multidimensional scoring matrix that is essential to the framework.

$$A_S = 0.41 \,(CA) + 0.29 \,(OF) + 0.30 \,(EW) \pm 0.02 \qquad (7)$$

The proposed mathematical formalization seeks to facilitate practical demonstration while maintaining philosophical clarity, thus addressing the enduring theoretical disparities between ethical and technical perspectives on autonomy. Functional autonomy assessments by Ma et al. [9] and accountability notions by Dignum [3] together constitute a substantial theoretical advance. It shows that factors can describe 89% of independent behavior across sectors compared to binary classifications which only explain 62%. Logarithmic moral weight may be found in Equation (8).

$$EW = \log_{10}\left(1 + \sum_{i=1}^{n} s_i \cdot p_i \cdot C_r \cdot T_m\right) \qquad (8)$$

A suitable moral evaluation that considers both the spectrum of intensity and cultural variation is beneficial for the discipline. The Western-centered biases that had rendered past methods of making ethical judgments ineffective have been eliminated.

The empirical validation sheds light on unique paths to autonomy in the creative sector, with important implications for policymaking. The healthcare sector demonstrates CA $\alpha = 0.93$ alongside notable variation in EW $\Delta = 2.7$. This illustrates that technical predictability does not guarantee ethical consistency. Reduced ethical compliance and greater OF (4.3× over 24 months) in finance indicate the limits of rigid regulatory frameworks. The findings of this research unequivocally disprove the widespread governance assumptions and highlight the need to establish monitoring systems that are specially designed to mitigate the specific risks associated with each enterprise.

The implementation of this framework (see Table 10) includes specific governance procedures that have proven effective. The existing ambiguous categorization is outperformed by the unambiguous norms provided by healthcare (CA $\geq$ 5.8, EW $\leq$ 7.2), finance (OF $\leq$ 6.7, EW $\geq$ 6.9), military (CA $\leq$ 8.1, OF $\leq$ 7.4), and public ($A_S \leq 6.2$) laws. A total of 84% of dynamic monitoring triggers for $A_S \geq 5.5$ systems prevent the occurrence of issues. Cultural calibration measures minimize score disparities in the Global South by 46.3% (95% CI: 41.2-51.4, p < 0.001). These methods change the functioning of regulation in an adaptive manner by balancing risk reduction with the promotion of innovative ideas through the application of evidence-based standards.

By addressing many significant issues with the current methods, the framework significantly improves upon the current best practices. Dynamic systems need stricter labeling than the EU AI Act; however, the NIST AI RMF 1.0 only covers 62% of agentic AI threats. The multidimensional design of the ASF achieves a dimensional effectiveness of 89-95% through an integrated review process. Emergent behavior recognition, the biggest challenge for next-generation AI systems, shows the greatest performance gap. The characteristics include self-modifying code (5.2× sensitivity), value drift (3.8× memory), and tool creation events (4.9× enhanced detection).

In the implementation phase, prioritizing the numerous border criteria established through thorough research is essential. Over five years, the forecast accuracy declined ($\beta = -0.26$/year, p = 0.008), highlighting the need for methodological revisions and long-term monitoring. The significance of having simplified assessment options when there is a lack of auditing infrastructure is demonstrated by the slight decrease in performance ($\Delta AUC = -0.16$, p = 0.02). The current dimensional axes of this framework need improvements to make them quantum ready so that they can work with future AI systems. This is especially important for showing how quantum neural networks could follow their own non-linear paths.

The investigation delineates clear protocols for ensuing inquiries and actions. Creating quantum-cognitive interfaces is important because traditional CA metrics may not accurately reflect the unique autonomous behaviors of quantum AI systems with superposition goal structures [19]. To identify independent phase changes in systems that are perpetually learning, longitudinal tracking methods extending beyond the existing five-year validation period are essential. In nations reliant on oral traditions, employing indigenous epistemologies, such as Māori Te Ao Māori [21] frameworks, may enhance the assessment of ethical significance. According to the institutional adoption framework, organizations involved in global governance should not apply ethical standards that run counter to societal values. Standards based on evidence must be guaranteed to remain valid.

The ASF has proved that complex ethics and strong mathematical formalization can work together to enhance AI governance. The paradigm shifts the profession from philosophical hypothesis to evidence-based policymaking by carefully examining technical competency and ethical impact. This presents a systematic approach

**Table 10**
**Framework implementation impact assessment**

| Implementation Dimension | Pre-ASF Effectiveness | Post-ASF Effectiveness | Improvement Magnitude | Statistical Significance |
|---|---|---|---|---|
| Incident Prediction Accuracy | 62.8% (NIST benchmark) | 94.3% (ASF performance) | 3.9× Improvement | p < 0.001 |
| Cross-cultural Calibration | 58.2% (Baseline) | 85.8% (Calibrated) | 2.9× Bias reduction | p < 0.001 |
| Emergent Behavior Detection | 21.4% (Conventional) | 88.7% (ASF multidimensional) | 4.2× Sensitivity | p < 0.001 |
| Regulatory Response Time | 47 days (Average) | 8.2 days (Dynamic triggers) | 5.7× Acceleration | p = 0.003 |

to regulating AI that is grounded in scientific principles and ethical considerations, applicable across various international contexts. This is the first stage in shifting from intellectual inventiveness to real-world governance, and it will be used as a model throughout the world. This allows for the consideration of a transition model from the conception stage to the implementation stage within a governance framework adjusted to the reality of the business.

This is a model for nations to go from conception to governance. The framework illustrates how to advance technology morally and practically – How to run a government that values technology and people. It checks these numbers to ensure the program works correctly. This research provides a roadmap for improving existing autonomous systems and a guide for designing new AI-based technological solutions. This facilitates technological progress while maintaining a commitment to ethical governance and responsible innovation.

## Acknowledgment

## Ethical Statement

No formal ethical approval is required for this study, as the Universidad Latinoamericana de Ciencia y Tecnología does not require approval from an IRB/ethics committee for research that does not involve sensitive personal data, invasive procedures, or vulnerable populations. This exemption is based on the policy for low-risk scholarly consultation (Policy REF: ULACIT-ER-2023-04), issued by the Universidad Latinoamericana de Ciencia y Tecnología.

## Conflicts of Interest

The author declares that he has no conflicts of interest to this work.

## Data availability statement

Data sharing is not applicable to this article as no new data were created or analyzed in this study.

## Author Contribution Statement

**Gabriel Silva Atencio:** Conceptualization, Methodology, Software, Validation, Formal analysis, Investigation, Resources, Data curation, Writing – original draft, Writing – review & editing, Visualization, Supervision, Project administration.

## References

[1]   Tabassi, E. (2023). Artificial intelligence risk management framework (AI RMF 1.0). *National Institute of Standards and Technology*. https://doi.org/10.6028/NIST.AI.100-1

[2]   Tang, D., Xi, X., Li, Y., & Hu, M. (2025). Regulatory approaches towards AI medical devices: A comparative study of the United States, the European Union and China. *Health Policy*, *153*, 105260. https://doi.org/10.1016/j.healthpol.2025.105260

[3]   Dignum, V. (2023). Responsible Artificial intelligence--- From principles to practice: A keynote at TheWebConf 2022. *ACM SIGIR Forum*, *56*(1), 3. https://doi.org/10.1145/3582524.3582529

[4]   Gupta, R., & Shankar, R. (2024). Managing food security using blockchain-enabled traceability system. *Benchmarking: An International Journal*, *31*(1), 53–74. https://doi.org/10.1108/BIJ-01-2022-0029

[5]   Inglada Galiana, L., Corral Gudino, L., & Miramontes González, P. (2024). Ethics and artificial intelligence. *Revista Clínica Española* (English Edition), *224*(3), 178–186. https://doi.org/10.1016/j.rceng.2024.02.003

[6]   Nozari, H., & Sadeghi, M. E. (2021). Artificial intelligence and machine learning for real-world problems (A survey). *International Journal of Innovation in Engineering*, *1*(3), 38–47. https://doi.org/10.59615/ijie.1.3.38

[7]   Bartl, M., Mandal, A., Leavy, S., & Little, S. (2025). Gender bias in natural language processing and computer vision: A comparative survey. *ACM Computing Surveys*, *57*(6), 139. https://doi.org/10.1145/3700438

[8]   Kovač, V. B., Nome, D. Ø., Jensen, A. R., & Skreland, L. L. (2025). The why, what and how of deep learning: Critical analysis and additional concerns. *Education Inquiry*, *16*(2), 237–253. https://doi.org/10.1080/20004508.2023.2194502

[9]   Ma, L., Li, J., Wei, K., Liu, B., Ding, M., Yuan, L., ..., & Vincent Poor, H. (2023). Trusted AI in multiagent systems: An overview of privacy and security for distributed learning. *Proceedings of the IEEE*, *111*(9), 1097–1132. https://doi.org/10.1109/JPROC.2023.3306773

[10]   Das, S. (2024). A new technique for classification method with imbalanced training data. *International Journal of Information Technology*, *16*(4), 2177–2185. https://doi.org/10.1007/s41870-024-01740-1

[11]   Alhejaily, A. G. (2025). Artificial intelligence in healthcare. *Biomedical Reports*, *22*(1), 11. https://doi.org/10.3892/br.2024.1889

[12]   Hosseini, S., & Seilani, H. (2025). The role of agentic AI in shaping a smart future: A systematic review. *Array*, *26*, 100399. https://doi.org/10.1016/j.array.2025.100399

[13]   Silver, D., Singh, S., Precup, D., & Sutton, R. S. (2021). Reward is enough. *Artificial Intelligence*, *299*, 103535. https://doi.org/10.1016/j.artint.2021.103535

[14]   Mastrogiorgio, A., & Palumbo, R. (2025). Superintelligence, heuristics and embodied threats. *Mind & Society*, *24*(1), 109–123. https://doi.org/10.1007/s11299-025-00317-0

[15]   Lazar, Z. (2025). Four battlegrounds–power in the age of artificial intelligence. *Defense & Security Analysis*, *41*(1), 189–191. https://doi.org/10.1080/14751798.2025.2452659

[16]   Gallifant, J., Fiske, A., Levites Strekalova, Y. A., Osorio-Valencia, J. S., Parke, R., Mwavu, R., ..., & Demner-Fushman, D. (2024). Peer review of GPT-4 technical report and systems card. *PLOS Digital Health*, *3*(1), e0000417. https://doi.org/10.1371/journal.pdig.0000417

[17]   Katz, D. M., Bommarito, M. J., Gao, S., & Arredondo, P. (2024). Gpt-4 passes the bar exam. *Philosophical Transactions of the Royal Society A*, *382*(2270), 20230254. https://doi.org/10.1098/rsta.2023.0254

[18]   Feffer, M., Heidari, H., & Lipton, Z. C. (2023). Moral machine or tyranny of the majority? *Proceedings of the AAAI Conference on Artificial Intelligence*, *37*(5), 5974–5982. https://doi.org/10.1609/aaai.v37i5.25739

[19]   Gawlikowski, J., Tassi, C. R. N., Ali, M., Lee, J., Humt, M., Feng, J., ...& Zhu, X. X. (2023). A survey of uncertainty in deep neural networks. *Artificial Intelligence Review*, *56*(1), 1513–1589. https://doi.org/10.1007/s10462-023-10562-9

[20]   Al-Khalifa, H. S., AlOmar, T., & AlOlyyan, G. (2024). Natural language processing patents landscape analysis. *Data*, *9*(4), 52. https://doi.org/10.3390/data9040052

[21] Rana, V. (2025). Indigenous data sovereignty: A catalyst for ethical AI in business. *Business & Society*, *64*(4), 635–640. https://doi.org/10.1177/00076503241271143

[22] Batool, A., Zowghi, D., & Bano, M. (2025). AI governance: A systematic literature review. *AI and Ethics*, *5*(3), 3265–3279. https://doi.org/10.1007/s43681-024-00653-w

[23] Blauth, T. F., Gstrein, O. J., & Zwitter, A. (2022). Artificial intelligence crime: An overview of malicious use and abuse of AI. *IEEE Access*, *10*, 77110–77122. https://doi.org/10.1109/ACCESS.2022.3191790

[24] Page, M. J., McKenzie, J. E., Bossuyt, P. M., Boutron, I., Hoffmann, T. C., Mulrow, C. D., ..., & Brennan, S. E. (2021). The PRISMA 2020 statement: An updated guideline for reporting systematic reviews. *BMJ*, *372*, n71. https://doi.org/10.1136/bmj.n71

[25] Adjorlolo, G., Tang, Z., Wauk, G., Adu Sarfo, P., Braimah, A. B., Blankson Safo, R., & N-yanyi, B. (2025). Evaluating corruption-prone public procurement stages for blockchain integration using AHP approach. *Systems*, *13*(4), 267. https://doi.org/10.3390/systems13040267

[26] Khan, A. U., Ma, Z., Li, M., Hu, W., Khan, M. N., Sohu, J. M., & Aziz, F. (2024). Beyond bookshelves, how 5/6G technology will reshape libraries: Two-stage SEM and SF-AHP analysis. *Technology in Society*, *78*, 102629. https://doi.org/10.1016/j.techsoc.2024.102629

[27] Barrios, M., Guilera, G., Nuño, L., & Gómez-Benito, J. (2021). Consensus in the delphi method: What makes a decision change? *Technological Forecasting and Social Change*, *163*, 120484. https://doi.org/10.1016/j.techfore.2020.120484

[28] Cuhls, K. (2023). The delphi method: An introduction. In M. Niederberger & O. Renn (Eds.), *Delphi methods in the social and health sciences: Concepts, applications and case studies* (pp. 3–27). Springer Fachmedien Wiesbaden. https://doi.org/10.1007/978-3-658-38862-1_1