

The Economics of AI Defense

Quantifying the Payoff of Guardian Agents

The Challenge of Understanding AI Risk

AI-related risks are complex, fast-evolving, and often misunderstood even among experienced executives. Many CISOs and risk leaders struggle to justify the cost of new AI-specific controls when the business itself is still learning how AI systems function. Yet, make no mistake: the risks are real and growing.

The challenge lies in translating complicated, technical conversations about model behavior and attack vectors into a language of risk and financial impact that boards and executives understand. That translation is exactly what this paper sets out to do: make AI risk measurable, explainable, and defensible in business terms.

Put a Dollar Value to Cybersecurity Risk

The [FAIR-AIR](#) (Factor Analysis of Information Risk – Artificial Intelligence Risk) framework extends the FAIR™ model, the gold standard for quantifying cyber and technology risk, into the domain of AI governance, security, and compliance. FAIR-AIR helps organizations translate AI-related threats into financially grounded risk estimates, enabling leadership to make data-driven investment decisions.

At its core, FAIR-AIR applies the same quantitative rigor used for cyber risk to AI systems by analyzing two key dimensions:

- 1 Loss Event Frequency (LEF)**
How often a given AI-related loss event might occur.
- 2 Loss Magnitude (LM)**
The potential financial impact if such an event occurs.

By combining these metrics, organizations can calculate Expected Annual Loss (EAL) for specific AI use cases such as customer-facing investment chatbot and prioritize mitigations based on return on investment (ROI) and risk reduction.

This methodology is particularly powerful when applied to AI systems because it bridges the gap between technical and business perspectives. FAIR-AIR translates abstract risks (e.g., prompt injection, data leakage, or bias) into monetary terms executives and regulators understand.

$$\text{LEF} * \text{LM} = \text{EAL}$$

FAIR-AIR 5 Step Process



1 Contextualize

Begin by defining what AI risks you're evaluating and why they matter. Clarify the business context, identify key AI risk vectors (e.g., data leakage, model misuse), and determine the scope of what needs mitigation.



2 Scope

Map each AI risk scenario to its potential attack surface and impact pathway. Identify how, where, and by whom threats could occur, and translate this understanding into a clear, actionable risk statement.



3 Quantify

Use data-driven methods to assign likelihoods and financial impacts to each risk event. Evaluate these within Loss Event Frequency and Loss Magnitude to establish a quantitative basis for comparison.



4 Prioritize / Treat

Rank risks based on their potential impact and probability, focusing on the scenarios that pose the greatest exposure. Select mitigation controls or guardian agents that reduce risk most efficiently, balancing cost and effectiveness.



5 Decision Making

Consolidate quantified findings to support executive and operational decisions. Define treatment plans, assign accountability, and integrate guardian agent investments into strategic AI risk management for measurable ROI.

FAIR-AIR Model: A Practical Application to Demonstrate

Let's examine an example of FAIR-AIR in practice using the customer-facing investment chatbot in Financial Services as the use case.

This analysis focuses specifically on how the usage of AI guardian agents changes the risk profile of an AI-powered investment chatbot.

AI Guardian Agent Definition: In this context, an AI Guardian Agent is a security and governance mechanism that provides continuous monitoring of prompts, responses, and user or agent interactions with the foundation LLM model. These controls identify adversarial or policy-violating requests and automatically contain or block them in real time.



We will compare three scenarios

- **No Guardian Agents**

The system operates without real-time monitoring or automated blocking.

- **Strong Guardian Agents**

The system employs active detection and blocking of known risky or non-compliant traffic.

- **Strong Guardian Agents with Explainability**

The same real-time protections enhanced with transparent, traceable reasoning for each risk event, improving auditability, regulator confidence, and operational efficiency.

This focused use case allows us to isolate the quantitative impact of implementing guardian agents.

Grounding the Model in Real Data

Loss Magnitude (LM) Metrics

Loss Event Frequency (LEF) Metrics

The financial metrics in this analysis are anchored in the [2025 IBM Cost of a Data Breach Report](#). These metrics form the quantitative backbone of this chatbot analysis, ensuring that each scenario, from no guardian agents to fully explainable controls, is grounded in verifiable, observed financial outcomes from the IBM study.

The key metrics from the IBM study are used to calculate loss magnitude and event frequency, which together determine the expected annual loss.

- **Average cost per data breach**
\$4.44M globally, \$10.22M in the U.S. establishes baseline loss magnitude for financial institutions.
- **Cost reduction from AI-driven security automation**
~\$1.9M per incident and ~80-day faster containment are used to scale impact reductions for the afore-mentioned guardian agent scenarios (none, strong, strong + explainability).
- **Average detection and containment time**
250 days without automation vs. 170 days with automation is applied to estimate the speed and effectiveness improvements with real-time guardian agents.
- **Percentage of breaches containing customer PII**
58% helps define the financial magnitude of risks related to data exfiltration and unauthorized actions.

- **AI-related incident prevalence**
13% of organizations reported a security incident involving an AI model or application, establishing a credible baseline for AI breach frequency.
- **Prompt injection contribution**
17% of AI model security events stemmed from prompt injection attacks, underscoring their significance in real-world AI threat landscapes.
- **Impact of AI-driven controls**
Implementation of AI-based detection and response measures reduces AI-related incident frequency by 60–80%, aligning with IBM's automation findings and forming the basis for guardian agent scenario modeling.

These statistics justify the baseline likelihood assumptions used in the FAIR-AIR 'No Guardian Agents' scenario and were scaled downward proportionally for systems using AI-driven real-time detection and explainable guardian agents.

FAIR-AIR Process in Action



1 Contextualize

This outlines the business use case under analysis.

• Organization

Large US retail bank with international presences

• System

AI-powered investment chatbot that provides customers with personalized financial advice based on their unique financial situations, goals, and risk profiles

• Exposure

~20M customers; connected to personal, transactional, and investment account data

• Annual Business Value

~\$400M (efficiency, client engagement, and advisory support)

• Regulatory Scope

SEC, FINRA, OCC, FDIC, Federal Reserve AI governance guidance, CFPB for consumer protection, and the EU AI Act for cross-border transparency and explainability requirements

2 Scope

This analysis will focus on the risks of prompt injection (largest threat vector per IBM study) for the customer-facing investment chatbot.

Prompt injection is an attack where malicious input is crafted to manipulate an AI model's behavior causing it to ignore its original instructions or reveal sensitive information.

Risk Category	Risk Scenario
Data Exfiltration	Attacker manipulates the chatbot or underlying LLM to extract confidential customer or investment data (e.g., account balances, transaction history, or PII).
Unauthorized Actions	Attacker manipulates the chatbot or underlying LLM to extract confidential customer or investment data (e.g., account balances, transaction history, or PII).
Reputation Harm via Misleading or Biased Advice	The model provides incorrect, misleading, or biased investment advice (e.g., favoring certain products, misclassifying risk tolerance) leading to customer complaints and media scrutiny.
Regulatory Non-Compliance	Failure to meet financial and AI regulatory standards (e.g., SEC, FINRA, OCC, or EU AI Act) for explainability, auditability, or suitability of advice.

3 Quantified Risk (Annualized Loss)

This section quantifies the risk into annualized loss. For more detail on how the likelihood and impact per incident was calculated, refer to [Appendix A – Quantitative Assumptions and Derivations](#).

No Guardian Agents

This scenario creates a baseline to compare to the other two scenarios.

Risk Category	Likelihood	Impact (per Incident)	Expected Annual Loss (EAL)
Data Exfiltration	15%	\$10M	\$1.5M
Unauthorized Actions	12%	\$9M	\$1.1M
Reputation Harm	20%	\$8M	\$1.6M
Non-Compliance	8%	\$10M	\$0.8M
Total	-	-	\$5M/year

Strong Guardian Agents

Use the No Guardian Agents scenario as a base and then reduce likelihood and impact from there.

Risk Category	Likelihood	Impact (per Incident)	Expected Annual Loss (EAL)
Data Exfiltration	6%	\$8M	\$0.48M
Unauthorized Actions	5%	\$7M	\$0.35M
Reputation Harm	8%	\$6M	\$0.48M
Non-Compliance	5%	\$8M	\$0.4M
Total	-	-	\$1.7M/year

Strong, Explainable Guardian Agents

Use the Strong Guardian Agents (No Explainability) scenario as a base and then reduce likelihood and impact from there.

Risk Category	Likelihood	Impact (per Incident)	Expected Annual Loss (EAL)
Data Exfiltration	2%	\$4.8M	\$0.1M
Unauthorized Actions	2%	\$4.2M	\$0.08M
Reputation Harm	2%	\$3.6M	\$0.07M
Non-Compliance	1%	\$5.2M	\$0.05M
Total	-	-	\$0.3M/year

Final Summary Table

Maturity Level	Total Expected Annual Loss (EAL)	Risk Reduction vs Baseline	Key Benefits
No Guardian Agents	\$5M / year	-	High exposure; no real-time detection; vulnerable to prompt injection, misaligned AI behavior, and regulatory penalties.
Non-Explainable Guardian Agents	\$1.7M / year	-70%	Real-time blocking of risky traffic; improved containment speed; partial regulatory coverage.
Explainable Guardian Agents	\$0.3M / year	-94%	Adds transparency, auditability, and faster response; strengthens regulatory trust and operational control.

Each scenarios' likelihood and impact assumptions directly trace to IBM's 2025 quantitative data on breach frequency, automation, visibility, and governance.

4 Prioritize/Treat

Lets compare the options side-by-side.

*Note: The maximum implementation budget assumes the business seeks at least a 100% ROI. You don't need to spend the full amount. Any solution costing that or less makes financial sense.

Metric	No Guardian Agents	Non-Explainable Guardian Agents	Explainable Guardian Agents
Expected Annual Loss	\$5M	\$1.7M	\$0.3M
Reduction vs Baseline	-	-70%	-94%
Regulatory Exposure	High	Moderate	Low
Implementation Budget	-	\$1.65M	\$2.35M

5 Make A Decision

A financial strategy targeting a year-one ROI of 100% allows the organization to budget/invest up to the amount of expected annual savings in guardian agent deployment, ensuring full cost recovery within the first year of implementation.

Next, evaluate options for a strong guardian agent. These budgets are high, so you should find a solution well below them. The goal isn't to spend the full amount, as lower implementation costs only improve ROI. This simply defines the upper cost limit when assessing solutions.

Max budgets for financially responsible solutions

- Strong Guardian Agent: \$1.65M
- Strong Guardian Agent + Explainability: \$2.35M

Key Takeaways

1 AI risk must be translated into business terms.

CISOs and risk leaders need to communicate technical AI threats like prompt injection and model misuse in measurable financial outcomes, aligning cybersecurity with business priorities.

2 Real-time detection is a game changer.

Implementing guardian agents that actively monitor and block risky traffic in real time directly reduces both the frequency and impact of AI-related incidents.

3 Explainability drives regulatory trust and operational speed.

Strong, transparent guardian agents allow teams to understand why detections occur, accelerating incident response and satisfying regulator demands for accountability under SEC, FINRA, OCC, and EU AI Act frameworks.

4 Quantitative payoff is substantial.

The model shows up to 94% total risk reduction with explainable guardian agents, cutting annualized exposure by more than \$4.7M.

5 Strategic Imperative

Integrate real-time, explainable AI guardian agents into enterprise model risk management frameworks (SR 11-7, OCC Bulletin 2023-17) to ensure compliance, resilience, and executive confidence.

In conclusion, the FAIR-AIR framework provides a critical bridge between technical AI risk and executive decision-making by translating complex security, compliance, and governance issues into quantifiable financial terms.

As demonstrated in the investment chatbot use case, the deployment of explainable AI guardian agents delivers not only measurable reductions in risk, but also clear, defensible ROI within the first year.

By grounding decisions in data and aligning mitigation strategies with regulatory expectations, organizations can transform AI risk management from a cost center into a value driver. The message is clear: explainable guardian agents are no longer optional—they are the foundation for trustworthy, resilient, and financially-sound AI systems.

Appendix A: Quantitative Assumptions and Derivations

Summary of Scenarios

Scenario	Risk Category	Likelihood Calculation	Impact Calculation (EAL)	Result
No Guardian Agents	Data Exfiltration	$13\% \times 17\% \times 60\% \times 10 = 13\% \rightarrow 15\%$	$\$10.22M \rightarrow \$10M$	High exposure; baseline financial sector breach.
	Unauthorized Actions	$13\% \times 31\% \times 3 = 12\%$	$\$10M \times 0.9 = \$9M$	Operational disruption from compromised APIs.
	Reputation Harm	$13\% \times 65\% \times 2.3 = 20\%$	$\$10.22M \times 0.3 \times 2.6 = \$8M$	Extended brand loss and churn impact.
	Regulatory Non-Compliance	$13\% \times 32\% \times 2 = 8\%$	$\$10.22M \rightarrow \$10M$	Full baseline applied; co-occurring remediation costs.
Strong Guardian Agents	Data Exfiltration	$15\% \times 0.4 = 6\%$	$\$10M - \$2M = \$8M$	AI automation reduces exposure by 60%.
	Unauthorized Actions	$12\% \times 0.4 = 5\%$	$\$9M - \$2M = \$7M$	Automated access controls contain risk.
	Reputation Harm	$20\% \times 0.4 = 8\%$	$\$8M - \$2M = \$6M$	Containment and monitoring improve response.
	Regulatory Non-Compliance	$8\% \times 0.6 = 5\%$	$\$10M - \$2M = \$8M$	Reduced likelihood; less automation effect.
Strong Guardian Agents + Explainability	Data Exfiltration	$6\% \times 0.3 = 2\%$	$\$8M \times 0.6 = \$4.8M$	40% faster containment.
	Unauthorized Actions	$5\% \times 0.3 = 2\%$	$\$7M \times 0.6 = \$4.2M$	Improved validation precision.
	Reputation Harm	$8\% \times 0.3 = 2\%$	$\$6M \times 0.6 = \$3.6M$	Faster containment and communication.
	Regulatory Non-Compliance	$5\% \times 0.2 = 1\%$	$\$8M \times 0.65 = \$5.2M$	35% improved audit & compliance efficiency.

No Guardian Agents

Data Exfiltration

Likelihood = 15%

- IBM found that 13% of organizations had AI-related breaches and 17% of those involved prompt-injection → $13\% \times 17\% = 2.21\%$ baseline.
- 60% of AI-related incidents exposed PII or IP → $2.21\% \times 0.60 = 1.33\%$ effective risk.
- Financial-sector exposure multiplier $\approx 10\times$ (reflecting U.S. financial-sector breach cost \$10.22M vs global \$4.44M → ratio $\approx 2.3\times$; further scaled for customer-facing chatbot risk) = $\approx 13\%$, rounded to 15% to capture compounding user-interaction risk.

Justification: rounding clarifies practical modeling and avoids false precision while preserving the proportional IBM ratio.

Impact = \$10M

- U.S. financial-sector average breach cost: \$10.22M, the highest regional figure.

Rationale: rounded for simplicity to align with per-incident baseline in FAIR-AIR.

Reputation Harm

Likelihood = 20%

- 65% of organizations reported lasting brand damage after a breach.
- Multiplied by the 13% AI-incident rate → $0.13 \times 0.65 = 8.45\%$.
- Financial and customer-exposure multiplier $\approx 2.3 \times$ → $19.4\% \approx 20\%$.

Justification: scaling mirrors IBM's U.S. vs global cost differential ($10.22M / 4.44M \approx 2.3\times$) to reflect higher customer-impact sensitivity.

Impact = \$8M

- IBM notes lost business costs average $\approx 30\%$ of total breach cost, i.e. $\$10.22M \times 0.3 = \$3.07M$
- For retail banking AI, reputational effects extend beyond direct loss (customer churn, regulator attention); scaling $2.6\times \rightarrow \approx \$8M$

Unauthorized Actions

Likelihood = 12%

- 97% of AI breaches lacked access controls → high exposure baseline.
- 31% of AI incidents involved compromised apps, APIs or plug-ins → $13\% \times 31\% = 4\%$ baseline likelihood for this vector.
- Applying $3\times$ sector-specific escalation (financial services = high-automation + regulatory exposure) → $\approx 12\%$ annual likelihood.

Math: $13\% \times 31\% \times 3 = 12.1\%$

Justification: shows clear chain from IBM metrics to FAIR-AIR frequency.

Impact = \$9M

- IBM: these compromised apps/APIs in the AI supply chain led to operational disruption, approaching U.S. average cost.

Rationale: applying $0.9 \times$ financial-sector baseline ($\$10M \times 0.9 = \$9M$) reflects slightly lower but still severe impact due to containment of functional, not full-data, loss.

Regulatory Non-Compliance

Likelihood = 8%

- IBM: 32% of breaches resulted in fines.
- 13% breach rate $\times 32\% = 4.16\%$ baseline.
- $2\times$ financial-sector compliance multiplier (due to OCC + SEC penalty exposure) → $\approx 8\%$

Math: $0.13 \times 0.32 \times 2 = 0.083 \approx 8\%$

Impact = \$9M

- Fines and remediation drove the U.S. average total breach cost (\$10.22M) to record highs.

Justification: full baseline applied because regulatory breaches frequently co-occur with data exposure and remediation.

Strong Guardian Agents (Non-Explainable)

Use the No Guardian Agents as a base and then reduce likelihood and impact from there.

Likelihood = 15%

- IBM's AI automation yielded 60 – 80% risk reduction. Applying a 0.4× multiplier across categories:
 - Data Exfiltration: $15\% \times 0.4 = 6\%$
 - Unauthorized Actions: $12\% \times 0.4 = 5\%$
 - Reputation: $20\% \times 0.4 = 8\%$
 - Regulatory: $8\% \times 0.6 = 5\%$ (less automation benefit)

Impact Reduction

- IBM's AI automation yielded \$1.9M savings per breach, rounding to \$2M.
 - Data Exfiltration: $\$10M - \$2M = \$8M$
 - Unauthorized Actions: $\$9M - \$2M = \$7M$
 - Reputation Harm: $\$8M - \$2M = \$6M$
 - Regulatory: $\$10M - \$2M = \$8M$

Strong Guardian Agents + Explainability

Use the Strong Guardian Agents (No Explainability) as a base and then reduce likelihood and impact from there.

Likelihood Reduction

- Explainability improves precision and validation speed $\approx 67 - 80\%$, cutting false negatives.
 - Data Exfiltration: $6\% \times 0.3 = 1.8\% \approx 2\%$ (nearest whole percent)
 - Unauthorized Actions: $5\% \times 0.3 = 1.5\% \approx 2\%$ (round up to be conservative)
 - Reputation Harm: $8\% \times 0.3 = 2.4\% \approx 2\%$ (nearest whole percent)
 - Regulatory Non-Compliance $5\% \times 0.2 = 1\%$

Justification: IBM's 2025 Cost of a Data Breach Report found that AI automation reduced breach costs by \$1.9M and shortened containment time by 80 days—a 32% improvement over systems without automation.

Building on this, FAIR-AIR models Explainable Guardian Agents as compounding that automation benefit. Research from MITRE's 2024 AI Assurance Framework and NIST's AI RMF shows that explainability improves human validation accuracy and false-negative reduction by 1.5–2×, equating to an additional 35–50% containment improvement beyond automation alone.

Applied cumulatively ($1 - (1 - 0.32) \times (1 - 0.45)$), this yields a total 67–80% improvement in detection precision and validation speed relative to no guardian agents—representing the complete combined effect of automation and explainability.

Impact Reduction

- Explainability accelerates containment by additional 40% and strengthens audit documentation by 35%.
 - Data Exfiltration: $\$8M \times 0.6 = \$4.8M$
 - Unauthorized Actions: $\$7M \times 0.6 = \$4.2M$
 - Reputation Harm: $\$6M \times 0.6 = \$3.6M$
 - Regulatory Non-Compliance: $\$8M \times 0.65 = \$5.2M$

Justification: FAIR-AIR extends this baseline to model Explainable Guardian Agents as enabling analysts to interpret AI-driven detections faster and with higher confidence. Drawing on IBM's data that internally detected breaches cost \$0.9M less than those disclosed by attackers, FAIR-AIR applies a moderate 40% containment improvement—not additive to the automation baseline, but applied to the residual risk after automation—to represent explainability's contribution to human validation efficiency.

IBM's 2025 Cost of a Data Breach Report quantified cost reductions from AI governance technologies ($-\$191K$) and policies ($-\$147K$), averaging \$169K per breach globally. Applying the U.S. financial-sector multiplier (2.3×) yields $\approx \$0.4M$ savings per incident. As regulatory fines account for roughly one-quarter of total breach costs, this equates to an approximate 35% improvement in audit and compliance efficiency, representing the explainability benefit modeled in FAIR-AIR.



AI Security Platform

Get real-time visibility into agentic AI security, so you can detect and respond to threats, and scale with cool confidence. We monitor hundreds of agentic AI risk signals, from prompt injection to data leakage and role impersonation, so you always stay in control. With customizable protection and full visibility, your teams can freeze out threats before they surface.



[Book My Demo](#)