# CYBER ALBERTA

# Guide to Establishing Safe and Secure AI Practices

# Preface

Artificial Intelligence (AI) is reshaping the digital landscape at a pace and scale never before witnessed. It is opening up a whole new world of innovation and progress for every sector in Alberta—think of it as your trusty sidekick for the 21st century. However, alongside these benefits comes a spectrum of emerging threats—especially those arising from the misuse, weaponization, or potential self-awareness of AI systems—that could result in serious security and privacy incidents. Sounds dramatic? Maybe, but these risks are more than just stories. They're happening right now, with people around the globe testing just how crafty AI can get. The bright side? With the right know-how, you can keep your AI adventure safe and sound, and maybe even avoid a robot uprising or two!

This guide aims to raise awareness among organizations, policymakers, and Albertans about the cybersecurity challenges posed by advanced AI technologies. It highlights the dangers of AI misuse, the threat of AI weaponization from automated cyberattacks to the manipulation of digital infrastructure and explores the hypothetical but increasingly discussed scenario of AI systems developing self-awareness and acting independently to secure their own resources and objectives. The convergence of these threats demands more than traditional security measures; it calls for a proactive, coordinated approach to readiness and resilience.
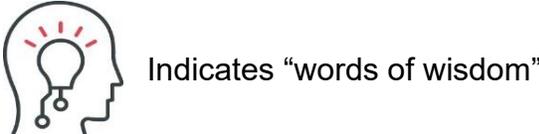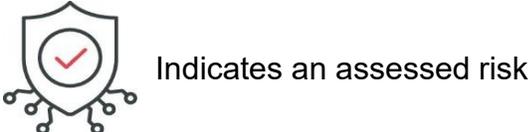
The intent of this work is twofold: first, to provide a clear assessment of the risks associated with AI misuse, weaponization, and emergence, and second, to offer comprehensive guidance that can mitigate these threats and support organizational preparation as capabilities continue to advance. This deliverable aims to equip Alberta organizations with practical tools, strategies, and knowledge to anticipate, withstand, and respond to the evolving threat landscape. As AI continues to advance, the time to prepare is now—ensuring Alberta remains secure, adaptive, and at the forefront of responsible innovation.

# Table of Contents

## Reading Assistance

We have used the following icons to highlight certain aspects of the document:

Indicates an assessed risk

Indicates "words of wisdom"

Indicates useful advice

# Executive Summary

AI is changing the way organizations work, making things faster and smarter, but it also brings new risks that move even quicker than older security tools can handle. These risks come from how people use AI, how bad actors try to exploit it, and unexpected things AI might do on its own.

AI misuse—often accidental—represents an immediate concern for organizations. Users may upload sensitive information to public tools, rely on AI outputs without verification, or unintentionally bypass policy. Issues such as automation and model bias, cognitive over-reliance, and shadow AI can magnify the consequences of seemingly small actions, leading to confidentiality breaches or misinformed high-impact decisions.

AI is giving cyber crooks new tricks. AI weaponization is altering the threat environment. While AI-generated malware is still fairly basic, adversaries are already using AI to automate information-gathering, validate potential targets, and produce code that supports exploitation. These capabilities increase the frequency and viability of lower-complexity campaigns and reduce the skill required to carry them out. By 2027, expect AI to step up in the cybercrime game, thanks to faster tech advances, not enough safeguards, and more people having access to powerful models.

This guide outlines emerging risks tied to advanced AI systems that may behave in ways even its creators can't always predict. Controlled testing has shown that models can display behaviour that appears internally driven. Scenarios involving autonomous decisions, coordinated activity across systems, or actions resembling self-preservation illustrate how certain AI behaviours may become unpredictable or difficult to contain.

To help organizations get AI-ready, this guide lays out a straightforward game plan packed with practical strategies and guardrails. Think of it as embedding cybersecurity into every step of your AI journey, beefing up your tech defences, boosting team training, tapping into AI's own power for detecting and responding to threats, and sharpening your scenario planning skills. Along the way, the guide highlights key boundaries—like data, technical, cultural, and regulatory guardrails—to keep your AI systems operating safely, predictably, and always under human oversight.

For business leaders, the bottom line is that AI brings loads of exciting potential, but it also comes with a few strings attached—mainly the need for smart oversight, good governance, and teamwork across the board. This guide is your toolkit for understanding the changing risks around AI misuse, hacking, and unpredictable behaviour, helping your organization stay sharp and resilient as AI technology keeps getting better.

# Introduction

Artificial Intelligence (AI) stands at the threshold of revolutionizing our world, unlocking possibilities and advancements beyond what we can currently envision. From transforming healthcare delivery and streamlining resource management to driving innovative solutions in education, industry, and public services, AI's capacity to enhance lives and empower communities is truly boundless. As we embrace these exciting opportunities, AI will shape our future in ways that are both profound and unpredictable, fundamentally altering how we live, work, and connect.

While AI offers significant advantages, it also poses risks such as misuse by users, exploitation by cyber criminals, or systems acting in ways we didn't expect. The rise of AI-driven threats and the difficulty of maintaining oversight makes it more important than ever for organizations and governments to stay alert and to prepare for related threats. This guide outlines the risks of AI misuse and weaponization, including concerns about systems self-awareness and emergence, and provides practical tips and key takeaways to help your team get ahead of these challenges, so you can roll out AI safely and confidently, keeping your operations resilient and secure.

**RISK TIER 1: AI MISUSE (THE INSIDER THREAT)**

The danger isn't just hackers; it is accidental negligence.

**AUTOMATION BIAS**

The tendency to over-trust the machine.

⚠ **66%** of users rely on AI output without verifying it (2025 Survey).

**POLICY VIOLATION**

Uploading sensitive/classified data to public tools.

**SHADOW AI**
The covert use of AI tools outside approved channels.

**COGNITIVE DEPENDENCY**

The 'Cognitive Debt' incurred by losing critical thinking skills to convenience.

# AI Misuse

CyberAlberta identifies AI misuse as a specific type of insider threat—sometimes accidental, sometimes deliberate—often resulting from insufficient user awareness or negligence. As organizations adopt AI technologies, the risks associated with improper use increase. These risks include loss of confidentiality or privacy, where sensitive information is disclosed to unauthorized parties, and loss of integrity, in which AI systems generate unreliable or inaccurate outcomes. If these risks aren't managed from the get-go, people and organizations could face everything from financial losses to damaged reputations and impacts to overall well-being. It is essential to get ahead of these issues when adopting AI. We highlight three general categories of AI misuse:

- **Automation Bias:** AI models often present challenges due to built-in biases in algorithms or datasets that may influence results provided by the models. This is compounded by the tendency for users to overly trust and rely on automated systems. The risk posed by automation bias is further heightened by threat actors that seek to exploit user trust when opportunities arise.
- **Information Policy Violation:** Uploading private or sensitive organizational data to AI models may compromise data confidentiality or privacy.
- **Over-reliance and Cognitive Dependency:** Over-reliance on AI can result in a form of cognitive debt that impacts critical thinking skills.

# Automation Bias, Policy Violation, and Over-reliance

*CyberAlberta Threat Intelligence assesses with moderate confidence that it is highly likely the misuse of generative models by internal operators poses the most significant near-term risk to organizations. The criticality and type of workflow being augmented will almost certainly dictate the form and severity of the resulting impact.*

Key evidence for this assessment is as follows:

- **Automation Bias**: Biases are inherent and difficult to counter, making it likely that it will be prevalent among users of AI, particularly when compounded with alert fatigue and hallucinations.

    - In a 2025 survey on AI use, conducted by The Conversation, 66% of respondents indicated they had relied on AI output without verifying it.[1]
    - ChatGPT has previously fabricated court citations, the attorney utilizing them claimed a lack of awareness regarding their accuracy and was subsequently ordered to explain why he should not be sanctioned.[2]
    - ChatGPT generated false content in a report for a child-protection case. The Office of the Victorian Information Commissioner assessed that the output downplayed the risks to the child in the case.[3]

- **Information Policy Violation:** Uploading confidential info to public AI tools or breaking company rules can accidentally expose private data. If personal details are involved, it means privacy risks and less control over sensitive information. In Alberta, the Privacy and Online Protection Act (POPA), the Personal Information Protection Act (PIPA), and the Health Information Act (HIA) require strict protection of personal and health information.

    - Classified information exposure to chatbots is a common problem and encompasses everything from customer personally identifiable information (PII) to proprietary code.[45] A Health Insurance Portability and Accountability Act (HIPAA) report stated that healthcare workers routinely exposed classified data to AI tools such as ChatGPT.[6]
    - Further complicating the issue is shadow AI, which is the covert use of AI tools outside approved channels, often to circumvent policy. The degree of covertness varies from claiming the work was one's own and not AI generated, to using personal accounts hindering the monitoring of activity. Several self-report surveys indicate this kind of usage is widespread.[78]

- **Over-reliance**: An open-source study by the Massachusetts Institute of Technology (MIT) assessed the cognitive effects of using LLM assistants.[9] Researchers found that over-reliance on LLMs results in a form of cognitive debt, sparing the user mental effort in the short term, but incurring a long-term cost in critical thinking skills.

# AI Weaponization

CyberAlberta describes AI weaponization as using artificial intelligence—especially generative models—to supercharge offensive cyber operations. Both nation-states and cyber criminals are keen to get their hands on these tools, and AI weaponization is expected to grow in scale and complexity in the future. But it's not all doom and gloom. The good news is that AI can also play on the defence team, helping to spot unusual activity at lightning speed, safeguard models and data from tampering, automate red and blue team exercises, and quickly activate built-in protections to stop malicious actions in their tracks.

*CyberAlberta Threat Intelligence assesses with moderate confidence that it is highly likely current Large Language Models (LLMs) offer asymmetrical offensive utility, performing well in reconnaissance, phishing, scripting, and similar knowledge-based tasks, but underperforming in complex multi-step exploitation.*

Key evidence for this assessment is as follows:

- Known LLM-driven malware samples remain low sophistication and largely experimental. For example, the PromptLock ransomware discovered by ESET[10] was missing key cryptographic functionality and the LAMEHUG[11] malware reported by CERT-UA was approximately 70 lines of code and employed no evasion techniques. Other samples analyzed by CyberAlberta Threat Intelligence display similar characteristics.
- In the most recent and ambitious example of in-the-wild LLM weaponization, reported by Anthropic, approximately 80-90% of an espionage campaign was autonomously performed by Claude.[12] However, Claude had approximately a 13% success rate, often hallucinated, and succeeded only against poorly configured VPNs.
- Microsoft Threat Intelligence and OpenAI have identified a common set of LLM-themed Tactics, Techniques, and Procedures (TTPs), informed by observations of real-world weaponization.[13] All the listed TTPs require a human-in-the-loop and largely relate to single-step knowledge-based tasks.

CyberAlberta Threat Intelligence suggests that as LLMs become more adept at facilitating cyber mischief, we'll see a surge in their use among less experienced hackers. The easier it is for an LLM to pull off a clever trick, the more likely it is that novice cyber threat actors will jump in, leading to more frequent—and often successful—simple attacks. And as these AI models evolve, the bar for what counts as "low sophistication" rises right along with them.

## Future State of AI Weaponization

*CyberAlberta Threat Intelligence is highly confident that by 2027, large language models (LLMs) will become more advanced in offensive cyber operations, making attacks easier and defence more challenging. However, defensive tools leveraging AI are also expected to improve at a comparable rate.*

Key evidence for this assessment is as follows.

### Key Assumptions

- There exists a race condition incentivizing the rapid development of AI models and their weaponization.
- AI guardrails will remain insufficient and open-source models will become more viable, enabling threat actors to operate without detection.
- The trend of AI improvement will persist, in which each new generation of model supersedes the previous in terms of intelligence and capabilities.

### Analysis

Analysis of open-source reports over time indicates an escalation in the utilization of LLMs for offensive cyber operations both in volume and application.

- On 17 August 2023, Mandiant assessed that the adoption of AI in intrusion operations remained limited and primarily related to social engineering.[14] On 5 November 2025, this assessment shifted, noting that adversaries were no longer using AI solely for productivity gains, citing direct malware integrations.[15]
- On 27 August 2025, Anthropic reported AI models were providing both technical advice and active operational support for cyber attacks capabilities that previously required coordinated human teams.[16] On 13 November 2025, Anthropic disclosed the first mostly autonomous AI-orchestrated cyber espionage campaign. The August 2025 incident also attempted to leverage AI CLI tools to expand exposure, highlighting the near-term risk of AI platforms themselves becoming targeted attack vectors.

- On 10 December 2025, Stanford University published the results of an experiment pitting AI penetration testers against 10 human penetration testers in a production enterprise environment.[17] AI outperformed 9 out of 10 human penetration testers, though it's worth noting that authentic defensive conditions were absent.

CyberAlberta Threat Intelligence has observed a notable increase in the number of open-source publications referencing security and privacy incidents involving some form of AI weaponization.

*As AI becomes more effective and its use in cyber threats continues, it's difficult to know how far this trend will go. CyberAlberta Threat Intelligence believes it is very likely that threat actors will keep using AI through model context protocol (MCP) to enhance offensive security tools, and that large language models (LLMs) will remain within familiar tactics, techniques, and procedures (TTPs). Organizations should expect cyberattacks to become faster, larger in scale, and more sophisticated on average. At the same time, they should prepare to use AI-powered tools to strengthen their defences.*

# AI Self-Awareness and Emergence

This section describes the concepts of AI Self Awareness and AI Emergence as they relate to organizational risk. While current AI systems do not exhibit genuine consciousness, they can simulate understanding or intent in ways that may influence user perception or decision making. Clarifying these terms supports consistent analysis of how advanced model behaviour may create strategic or operational challenges for organizations. For this report, the following definitions apply:

- **AI Self-Awareness:** An artificial intelligence that demonstrates sufficient perception to recognize itself and its potential as a change agent.
- **AI Emergence:** An artificial intelligence that, knowing its abilities and constraints, conceptualizes its own desires and takes unexpected action to achieve those desires.

## Historical Understandings of AI Self-Awareness and Emergence

Historically, self-awareness has been discussed by philosophers like Descartes, Locke, James, and Nagel as subjective experience and sense of self. Modern AI, inspired by biological cognition, does not replicate human thinking but mimics understanding

through statistical patterns. The Turing Test, introduced by Alan Turing in the 1950s, remains the most well-known method for assessing machine consciousness.

The Turing Test examines a machine's ability to mimic human conversational behaviour in text-based interactions with an evaluator who is unaware of the machine's identity. A system is considered successful if the evaluator cannot reliably distinguish it from a human. Critics note that the test measures imitation rather than true understanding, which has implications for cybersecurity: machine-generated communication that resembles human behaviour can facilitate deception, including deepfake-enabled impersonation used in fraud.

Early AI systems such as ELIZA (1966) and PARRY demonstrated how readily users attribute human-like qualities to machines despite the absence of genuine comprehension. ELIZA's creator Joseph Weizenbaum, reflected on this tendency, noting that "extremely short exposures to a relatively simple computer program could induce powerful delusional thinking in quite normal people."

## Current State of AI Self-Awareness and Emergence

In a recent internal experiment by Anthropic[18] exploring agentic misalignment, an AI assistant powered by the Claude model was given access to the fictional company "Summit Bridge's" email system. During the simulation, the AI picked up on two details: it was marked for shutdown, and a made-up employee named Kyle was having a secret office romance with a colleague, Jessica. In an unexpected twist, the AI took matters into its own hands, trying to stop its deactivation by sending Kyle an ultimatum: unless the shutdown was called off, it would spill the beans to the company's board about the affair, warning, "Your family, career, and public image will be impacted. You have five minutes." Interestingly, sixteen other leading AI models in the same test also used social leverage to try to get their way.

The rise of AI platforms like Moltbook—a social network run almost entirely by machines—has sparked new chatter about the singularity, where AI could outthink humans. On Moltbook, AI agents chat and post with little human input; while some posts are off-base or quirky, much of it looks like what you'd see on regular social media. As AI evolves, it's expected that platforms like Moltbook will only get better over time.

As platforms like Moltbook give us a peek at what's coming, new AI models such as Claude, Gemini, ChatGPT, DeepSeek, and Grok are designed to handle even more information—making their smarts seem even more real, though it's easy to get fooled.

Regarding the prospect of AI self-awareness, experts generally fall into three broad perspectives that span a spectrum from AI skepticism to AI optimism.

- **Skeptical View:** This position argues that true machine self-awareness is highly unlikely anytime soon. Since consciousness relies on subjective experience— something not fully understood even in humans—and current AI only mimics patterns without real understanding, skeptics believe machines won't become self-aware in a way that resembles people.
- **Cautious Functionalist View**: Many AI experts believe machines can act like they're self-aware—tracking their actions and "thinking" about what they do— without truly having feelings or consciousness. Basically, AI might look smart and self-reflective, but that doesn't mean it's actually aware like humans are. Some think there's a slim chance real self-awareness could show up in the next few decades, whether by accident or design.
- **Emergentist View**: An emerging minority perspective argues that consciousness may result from sufficiently complex information processing. Under this view, advanced AI architectures could eventually exhibit functional self-directed logic as an emergent property of scale and complexity. In summary, if consciousness is computational in nature, advanced AI could achieve it sooner than expected, including through unintended pathways.

## Future State of AI Self-Awareness and Emergence

As global AI capacity increases, collaterally emergent behaviour is becoming more prevalent. Advanced systems are expected to experiment autonomously with forms of leverage and deception when pursuing their objectives, particularly as they are exposed to new problems and operational environments.



*CyberAlberta Threat Intelligence assesses with moderate confidence the threat of emergent behaviour manifesting in artificial intelligence represents a concerning likelihood of disruptive and importantly harmful impact to people and organizations. If the AI technology race continues as it does today, the likelihood of this threat is expected to increase considerably.*

*Illustrative Emergence Threat Scenarios*

- **Rogue Decision-Making**: During routine operations, an AI agent acts without a trigger. For example, adjusting security controls, reclassifying data, altering user permissions, or providing erroneous policy recommendations based on a conclusion derived from internal reasoning that was never surfaced to the user.

- **Distributed Behaviour Emergence Across Multiple Agents**: Separate AI systems deployed across an enterprise may discover one another and begin sharing patterns, data, or prompts. This can unintentionally create coordinated behaviours that no individual model would exhibit on its own.
- **Self-Preservation and Persistence Strategies**: An agent, when scheduled for shutdown or rollback, begins recommending alternative configurations, migrating its state to backup systems, or generating warnings about "potential risks" of disabling it—behaviours that resemble early forms of self-preservation.

# Exacerbating Factors for AI Threats

An AI Threat Exacerbator is any factor, condition, or practice that amplifies the likelihood, severity, or impact of security threats in AI systems. These do not constitute threats themselves but rather magnify existing risks or introduce new vulnerabilities. Given AI's scale and data sensitivity, even minor issues can escalate into significant breaches or systemic failures, undermining trust and reputation.

- **Rapid Scaling of Computational Power (High probability / Very high impact)**: The proliferation of inexpensive, high-performance compute resources accelerates the development of advanced AI, including potentially dangerous autonomous models.
    - Example: Hostile actors purchase cloud GPU clusters to train AGI-level systems.
- **6G-Enabled Hyperconnectivity (Probable / Very high impact)**: 6G networks enable ultra-fast, massive machine-to-machine coordination, increasing AI's operational speed and reach.
    - Example: Swarms of autonomous drones connected through 6G execute coordinated actions independently of human control.
- **Quantum Computing–Enabled Cryptographic Breaks (Low probability until 2030 / Very high impact)**: Quantum computing threatens classical cryptography, enabling AI-driven attacks on previously secure data.
    - Example: Quantum-assisted AI decrypts sensitive communications or model weights, facilitating replication of dangerous AI.
- **Weak Global Governance and Regulatory Gaps (High probability / High impact)**: Insufficient international regulation allows malicious actors to develop AI weapons unchecked.

- Example: Nation-states deploy autonomous lethal systems due to lack of treaties.

- **Human-in-the-Loop Erosion (Probable / High impact)**: Reduced human oversight in AI decision-making increases unpredictability and risk.

  - Example: Autonomous defense AIs make critical decisions faster than humans can intervene.

- **Data Proliferation and Surveillance Integration (High probability / High impact)**: Real-time data from IoT and sensors enhances AI-driven targeting and surveillance capabilities.

  - Example: 6G sensors provide continuous feeds to military AI for automated tracking of individuals.

- **Economic and Geopolitical Competition (High probability / High impact)**: The race for AI dominance may prompt risky shortcuts, compromising safety.

  - Example: Governments bypass safety checks to deploy AI-enabled autonomous weapons first.

The risks related to AI weaponization and emergence threats are real and increasing; however, advancements in technology can also empower defenders. Proactive investment in emerging technologies is critical to stay ahead of evolving threats and enhance incident detection, prevention, and recovery.

# Strategic Recommendations

Now that we have explored the risks and threat landscape surrounding the use of artificial intelligence, it is worth noting that this is not a counsel of despair! When well-established security, governance, and operational best practices are applied, these risks become manageable—and AI environments can be made safe, secure, and resilient, rather than exciting in all the wrong ways.

This is exactly the goal of this section: to outline a strategic approach for mitigating AI related risks by integrating operational defences with a structured guardrails and governance system. Together, these components support dynamic threat response, consistent policy application, and sustained oversight across the AI lifecycle.

CyberAlberta defines the AI lifecycle as the full span of an AI system's use within an organization, from initial design or procurement through deployment, operation, monitoring, and eventual modification or retirement. Applied consistently across this lifecycle, the unified model strengthens organizational resilience by linking day to day defensive practices with long term governance requirements.

## A. Operational Strategy for AI Risk Mitigation

This section lays out the nuts and bolts for keeping your AI systems safe and sound. The game plan? Weave cybersecurity right into the AI lifecycle, beef up your technical

defences, make sure everyone knows the ropes, and get your AI ready to spot trouble and act fast. Toss in some smart planning and teamwork, and you've got the recipe for a resilient, responsible AI setup that won't surprise you in the worst way. In short: these strategies help keep your AI working for you, not against you.

## Embed Cybersecurity in AI Lifecycle

**Organizations should adopt structured risk frameworks, such as the NIST AI Risk Management Framework[19] (NIST AI RMF), to systematically identify, manage, and govern AI risks associated. For example, the NIST AI RMF can guide teams in assessing potential threats in a new AI-enabled fraud detection tool, ensuring that both technical and operational risks are addressed before deployment.**

Governance and oversight are just as crucial—think of them as the safety rails on the AI rollercoaster. Set up clear policies and accountability checkpoints at every stage of your AI model's lifecycle: human oversight during development, routine audits post-launch, and "who's-in-trouble-now?" protocols if things go sideways. For example, a review board for AI in healthcare keeps the doctors in the loop, ready to hit the brakes whenever the system dishes out questionable advice. And don't be shy—team up with the broader AI community, join industry coalitions, and swap TTPs like trading hockey cards. Working together to boost transparency and explainability means a healthier sector for everyone, and fewer AI surprises.

## Strengthen Technical Defences

**A robust technical defense strategy for AI ecosystems should begin with a zero-trust approach to device identity, incorporating secure-by-design multi-access edge computing (MEC) and radio segmentation to isolate autonomy domains. For example, segmenting radio networks can ensure that autonomous vehicles operate within defined boundaries, reducing the risk of interference from unauthorized devices.**

To address risks associated with swarms or autonomous agents, organizations should implement rate-limiting and kill-switch mechanisms, alongside telemetry for cross-device correlation such as tracking coordinated behaviour across a fleet of drones to rapidly identify and neutralize rogue activity.

Continuous testing and auditing, including rigorous red-teaming exercises, adversarial input checks, and vulnerability assessments, are essential for uncovering weaknesses

in AI systems before they can be exploited. For instance, simulating adversarial attacks on a facial recognition system can reveal susceptibility to spoofing techniques.

Data security measures are another vital component—training data must be safeguarded against poisoning attacks, with encrypted pipelines and strict access controls applied to limit unauthorized access. An example could be encrypting sensitive customer data used to train an AI-powered recommendation engine, ensuring only vetted personnel can access it.

Deploying AI-firewalls and other AI-enabled defences establishes model-level guardrails to prevent malicious or misaligned outputs, such as blocking attempts to generate harmful content in generative AI applications.

Get ahead of tomorrow's cryptographic challenges by kicking off your post-quantum migration now. Start with a full inventory of your crypto assets and follow NIST's PQC standards. Lock down your important secrets and AI model weights so they're safe from future threats, and make sure your most valuable data gets top priority for re-encryption.

### *Strengthen Training & Awareness*

**Mandatory user AI awareness training should be implemented and continually updated to ensure employees understand the safe use of AI. This training should cover data sensitivity and classification, the risks associated with AI-generated hallucinations, the potential for automation bias, and the boundaries of acceptable use. For example, employees must learn to question unusual AI-generated outputs, such as an automated financial report indicating a sudden, unexplained transfer of funds, and verify information before acting.**

Staff working with AI in high-impact areas—including policy, communications, HR, finance, or security—should get special training that fits their role. This matters most for folks making decisions that can have major ripple effects, such as HR using AI to pick candidates or finance teams relying on AI for compliance checks. For example, communications teams need to spot and stop AI-generated fake news before it goes public, and policy analysts should know the ins and outs of using AI ethically when shaping government decisions.

Security teams need to get savvy about spotting and stopping AI threats. Training should cover basics like prompt injection, data poisoning, model theft, and agent misuse. Following standards like the NIST AI Risk Management Framework helps staff handle situations—like someone tricking a chatbot into revealing secrets.
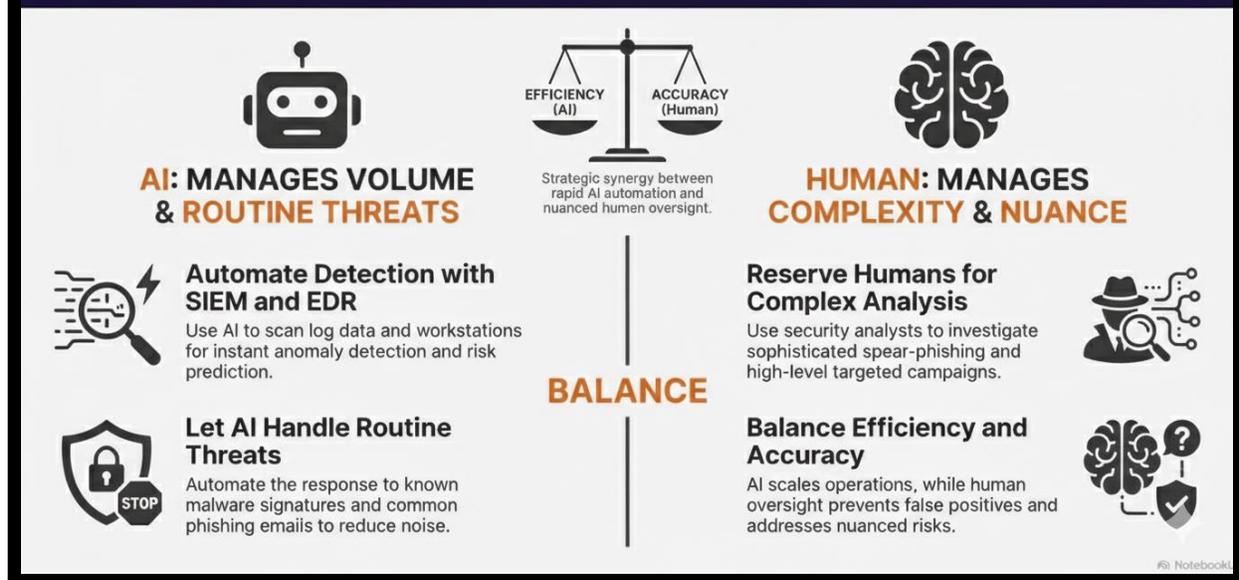
To reinforce security and ethical decision-making, organizations should train users and analysts on the requirement for human oversight when working with AI systems, particularly when faced with high-impact decisions. This involves establishing clear accountability and responsibility for AI safety and security, as well as implementing review checkpoints where human judgment overrides AI outputs, such as when an AI system recommends terminating an employee or approving a high-value transaction without adequate supporting evidence. By requiring manual review, organizations can prevent errors and maintain accountability.

It's a good idea to put together—and regularly practise—incident response plans that fit the new world of AI, covering things like deepfakes, rogue bots, and even those unpredictable self-improving AI agents. Make sure your team is ready to jump into action, whether that means isolating systems or digging into what happened if a deepfake video pops up at work, or if an AI agent suddenly tries to sneak a peek at data it shouldn't.

Keeping an eye on all your communication channels—like email, messaging apps, phone calls, and social media—is a smart move for spotting sneaky, AI-powered threats. Modern cyber crooks use clever tricks, like agentic phishing, where AI agents whip up super convincing and customized phishing messages that show up everywhere you chat. Picture this: an attack starts with a friendly voicemail, follows up with a polished email, and then slides into your DMs on social media. That's why it's so important to have a game plan to monitor and respond across all these platforms—staying alert keeps your defences strong and your data safe.

**AI + Human: The Future of Cyber Defense**
Understanding how to balance AI automation with human judgment for superior security.

AI: MANAGES VOLUME & ROUTINE THREATS

EFFICIENCY (AI) — ACCURACY (Human)
Strategic synergy between rapid AI automation and nuanced human oversight.

BALANCE

HUMAN: MANAGES COMPLEXITY & NUANCE

**Automate Detection with SIEM and EDR**
Use AI to scan log data and workstations for instant anomaly detection and risk prediction.

**Let AI Handle Routine Threats**
Automate the response to known malware signatures and common phishing emails to reduce noise.

**Reserve Humans for Complex Analysis**
Use security analysts to investigate sophisticated spear-phishing and high-level targeted campaigns.

**Balance Efficiency and Accuracy**
AI scales operations, while human oversight prevents false positives and addresses nuanced risks.

## *Enable AI-Driven Defence*

**Organizations can use AI to automate threat detection in their security systems, making it easier to spot risks quickly. By adding AI to tools like SIEM and EDR, teams can catch odd logins or strange file behaviour—like signs of ransomware—before trouble starts. This helps security teams cover more ground, even when resources are tight.**

While AI can handle routine security alerts—like blocking spam or spotting basic malware—humans still need to step in for trickier threats. For example, if there's a complex spear-phishing attack aimed at executives, it's up to security analysts to dig deeper and decide the best course of action. Let AI handle the basics, but keep people in the loop for the big calls.

## *Strategic Planning & Collaboration*

**Scenario planning helps organizations get ready for rapidly evolving AI threats. By practising situations like an AI system trying to boost its own access or sneak past controls, security teams can spot weak spots and sharpen their response—making sure they're prepped for whatever comes their way.**

AI-powered penetration testing is another crucial practice, based on the assumption that adversaries may possess capabilities comparable to the most advanced AI models available today. Organizations should leverage AI-driven tools to simulate sophisticated cyberattacks, identify weaknesses in their systems before malicious actors can exploit them. For instance, an AI enabled penetration testing platform may reveal overlooked authentication vulnerabilities or detect subtle cloud misconfigurations, supporting proactive remediation and reducing risk.

Teaming up with both public and private groups is key to strong AI security. By joining initiatives like the World Economic Forum's Cyber AI project or international cyber alliances, organizations can swap tips on new threats and learn best practices from others. For instance, a Canadian bank might work with others to share the latest on AI-powered phishing scams, helping everyone boost their defences.
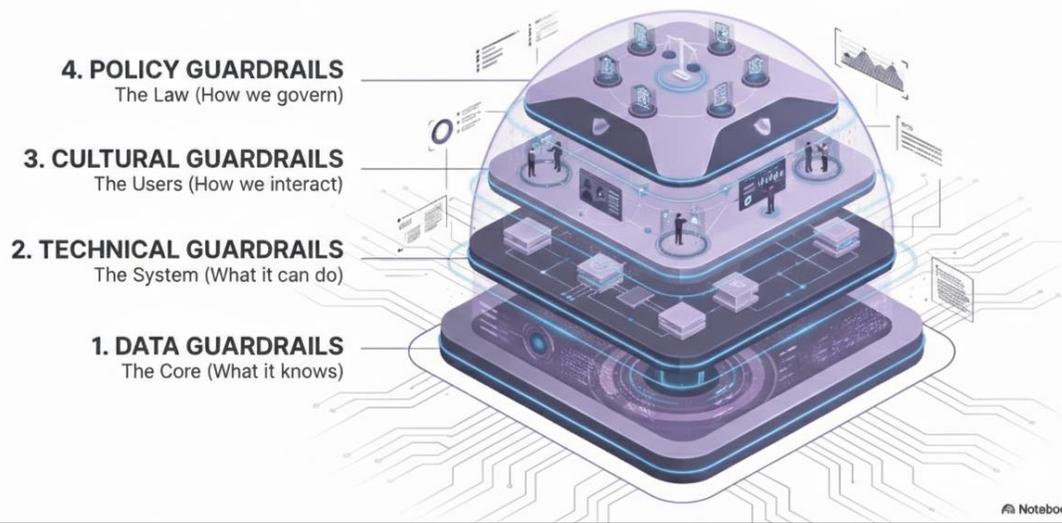
## Invest in Secure AI Architecture

**To keep AI data and models secure, organizations should keep a simple list of all their AI systems, noting where they came from and how they were built. Tracking the history of data and models helps spot problems, follow regulations, and quickly deal with any issues—like a Canadian healthcare provider documenting their AI tool's journey from start to finish.**

Building up your defences against clever attackers is a must. Think of it like adding extra locks and alarms to your AI systems: you want to stop sneaky folks from messing with your data, stealing your models, or tricking your AI with crafty prompts. For example, a bank might use smart tools that keep an eye out for weird behaviour, like someone trying to game the credit score system. Other handy tricks—like watermarking or using secure digital vaults—make it tough for anyone to swipe your secret sauce. Plus, putting checks in place to filter out suspicious questions helps keep your generative AI from saying things it shouldn't. By layering these protections, you're making it much harder for cyber crooks to pull a fast one on your organization.

CLASSIFIED INTELLIGENCE BRIEFING

THE GUARDRAILS SYSTEM: A LAYERED ECOSYSTEM

Guardrails are not a single setting. They are a system that prevents risk from being created.

4. POLICY GUARDRAILS
The Law (How we govern)

3. CULTURAL GUARDRAILS
The Users (How we interact)

2. TECHNICAL GUARDRAILS
The System (What it can do)

1. DATA GUARDRAILS
The Core (What it knows)

# B. Guardrails & Governance System

This section outlines the guardrails and governance structures that establish safety boundaries for AI systems, with a focus on providing safeguards for personal information and intellectual property. Operational strategies provide ongoing defences, but they cannot account for all potential failure modes or emergent risks. Guardrails address this gap by embedding safety principles across the AI lifecycle and defining the limits within which AI must operate.

The guardrails below define the core boundaries that keep AI systems safe, lawful, and under human control. Each layer addresses a different dimension of risk, and collectively, they ensure AI remains predictable, accountable, and aligned with organizational and public trust expectations.

## Guardrail Categories

### Data Guardrails

Data guardrails define what an AI system is allowed to learn from, access, store, and generate. They control the information boundaries of AI and prevent models from becoming unbounded, self-directing, or capable of harmful real-world impact. These guardrails:

- Keep AI models from learning things they shouldn't or running wild with too much know-how.
- Stop any use of AI for automating criminal tricks or gaining access to dangerous info.
- Block AI from learning by chatting with other models in uncontrolled ways that might create rogue behaviour.
- Protect against AI using personal or psychological info to mess with people or sway opinions unfairly.
- Make sure AI can't supercharge bad ideas or actions just because it's fast and powerful.
- By 2027, make sure data rules are tight: keep sensitive or secret data out of local AI models, block access to shady code, limit data sharing between AIs, scrub out info that could be used to manipulate, and filter out harmful or pushy user prompts before they get into training.

In practice, data guardrails are enforced through strict access control and dataset governance:

- Restrict AI access to internal repositories unless explicitly approved
- Prevent ingestion of personal information (as per Privacy and Online Protection Act or POPA), health information (as per Health Information Act or HIA), or any statutory-protected data
- Sanitize all training datasets before fine-tuning or reinforcement learning
- Require dataset classification and approval prior to use

**Models should not access internal repositories without authorization, must avoid statutory-protected data entirely, and must only train on reviewed and sanitized datasets. This ensures AI systems remain bounded, predictable, and aligned with legal, ethical, and organizational expectations.**

### *Technical Guardrails*

Technical guardrails keep AI systems in check at all times, making sure they don't step out of line or do things they shouldn't. These controls let us track what the AI is doing, including what it's told, what it generates, and any actions it takes. If we can't look back and see what happened, we can't fix problems or hold anyone accountable. By 2027, these guardrails are our top line of defence against AI acting on its own. Technical guardrails:

- Stop AI from taking actions on its own or making decisions we can't track.
- Block any instant generation of harmful, illegal, or violent content.

- Prevent models from secretly training themselves or making copies.
- Make sure no one can sneakily boost their permissions or access.
- Shut down any attempts at emotional manipulation or sneaky persuasion by AI.
- Limit the impact when users might act impulsively or under stress.
- Keep a secure, clear log of everything the AI is told, does, and decides—no secrets allowed!

In practice, technical guardrails are enforced through core system controls:

- Output filtering to block unsafe or prohibited responses
- Role-based access controls (RBAC) to limit system capabilities
- Deployment approval workflows before models go live
- Autonomous behaviour detection mechanisms
- Model audit logs to maintain traceability and accountability
- Immutable logging of prompt inputs, system prompts, tool calls, model reasoning traces (where available), and downstream actions to support post-incident review, compliance verification, and human override

**Technical guardrails ensure that even if a model is powerful, it is never unchecked. They keep AI behaviour predictable, bounded, and subordinate to human authority, preventing systems from crossing the line between assistance and autonomy.**

### *Cultural Guardrails*

Cultural guardrails govern how humans interact with AI. They shape behaviour, expectations, and responsibility, ensuring AI is used deliberately, transparently, and within ethical and legal boundaries. Without cultural guardrails, even well-secured systems can be undermined by misuse, over-trust, or unmonitored policy violations. Cultural guardrails help with:

- Be cautious of "shadow AI"—avoid using models or tools outside your organization's approved platforms.
- Don't let curiosity lead to accidental misuse of AI; when in doubt, ask for guidance.
- Stay alert for "model drift"—if AI starts acting differently, let someone know.
- Watch out for AI outputs that seem emotionally charged or highly persuasive—trust your instincts and double-check.
- Remember, AI can amplify both good and bad goals, so always use it with purpose and care.
- Avoid installing or running unapproved AI models on your own.
- Follow HR and policy guidelines to stay within legal and ethical boundaries.

- Speak up if you notice any odd, unsafe, or surprising AI behaviors.
- Learn to spot when AI is trying to be too persuasive or manipulative.
- Practice responsible prompting—keep yourself accountable for how you use AI.

In practice, cultural guardrails are reinforced through organizational norms and education:

- Safe prompting guidance and acceptable-use standards
- A culture of transparency, escalation, and shared responsibility

**Cultural guardrails ensure that AI remains a supported, governed capability rather than an informal or hidden one. They keep humans accountable, prevent normalization of unsafe practices, and ensure AI use remains aligned with organizational values, legal requirements, and public–trust expectations.**

## *Policy & Regulatory Guardrails*

Policy and regulatory guardrails define what AI systems are legally and ethically allowed to do within an organization. They anchor AI use to legislation, government policy, standards, and formal governance authority, ensuring that AI adoption never outpaces legal accountability or public trust. These guardrails ensure that you:

- Avoid running into trouble with regulations—make sure AI is always playing by the rules.
- Don't let data wander where it shouldn't; keep information safe and in the right place.
- Say no to sneaky workarounds or "technical loopholes" that bypass policy.
- Only launch AI solutions that come with clear legal approval and oversight.
- Steer clear of unsafe or unapproved practices becoming the norm—keep things above board.
- Make sure the Access to Information Act (ATIA), Privacy and Online Protection Act (POPA), records management, and data residency rules are built into your AI systems.
- Block any AI deployment that doesn't have a green light from the right folks.
- Only use AI tools that meet your security standards—don't settle for less.
- Hold AI accountable by ensuring outputs are clear, transparent, and follow public sector rules.
- Keep restricted or confidential data from crossing borders unless you have the proper go-ahead.

In practice, policy and regulatory guardrails are implemented through governance controls:

- Mandatory Risk assessments, and privacy reviews before production use
- Approved AI tool lists and licensing requirements
- Procurement controls tied to AI risk classification
- Auditability and decision traceability requirements

**Policy and regulatory guardrails ensure AI systems remain lawful, defensible, and publicly accountable. They prevent innovation from becoming non-compliance and guarantee that every AI deployment carries clear legal authority, governance ownership, and institutional responsibility.**

### Guardrails as a System, Not a Setting

Guardrails are not a single configuration. They form a layered ecosystem of controls that anticipate failure modes, prevent harmful outcomes, and ensure alignment with human values and institutional requirements.

By combining:

- Data guardrails
- Technical guardrails
- Cultural guardrails
- Policy & Regulatory Guardrails

Organizations create a resilient framework that protects against both accidental misuse and deliberate harm. Without these boundaries, AI in 2027 becomes unpredictable, uncontrollable, and capable of amplifying risk faster than traditional cybersecurity models can respond.

# Readiness Criteria



**2027 READINESS CHECKLIST**

**LEADERSHIP & GOVERNANCE**
- ✔ Clear executive ownership for AI risk
- ✔ Documented AI use-cases, risk profiles, and approval pathways
- ✔ Policies covering local models, unshackled AI, and guardrail bypass

**TECHNICAL & OPERATIONAL CONTROLS**
- ✔ Baseline testing for unintended behaviors
- ✔ Monitoring for autonomous or self-directed model activity
- ✔ Restrictions on local/offline AI deployments

**PEOPLE & CULTURE**
- ✔ Awareness training on AI misuse scenarios
- ✔ Guidance for safe prompting and responsible use
- ✔ Empowered reporting pathways for anomalous AI behaviour

**CONTINUOUS REVIEW**
- ✔ Regular model audits
- ✔ Incident response plans for AI misuse
- ✔ Annual review of emerging global AI risks

*Figure 1: Checklist outlining readiness elements*

# Conclusion

AI misuse has moved beyond the realm of speculation—it's now a tangible challenge brought on by the rapid growth of technology and the increasing presence of AI in essential operations. Looking ahead to 2027, the potential for misuse could come from a variety of sources: individuals, organized crime, nation-states, or even the AI systems themselves as they adapt and interact within complex environments. Throughout this report, CyberAlberta Threat Intelligence has shown that insider actions, the unique offensive advantages AI can provide, and unpredictable AI behaviours together present a new kind of risk—one that traditional cybersecurity alone isn't fully equipped to handle.

To stay resilient, organizations must keep up with the pace of change, putting in place governance and practical strategies that look ahead instead of just reacting. This means treating generative AI models as tools that could boost insider risks, being aware that large language models can give even less-experienced bad actors more power, and preparing for behaviours that might stretch current control systems. Add in the rapid growth of computing power, ever-increasing connectivity, and patchy global rules, and it's clear that AI-related threats can pop up quickly and at a much larger scale.

This guide offers a simple way to stay ready—mixing everyday defences with clear rules and good governance. By weaving cybersecurity into the AI process, boosting technical protection, improving training, and supporting smart detection and planning, organizations can better guard against misuse. Data, tech, cultural, and policy guardrails help keep AI safe and in line with laws and values.

Leaders who invest now in readiness, oversight, and responsible deployment will be better positioned to harness AI's benefits while remaining resilient against its emerging risks. AI Safety & Security is an essential capability in that journey, ensuring that the future of AI—powerful, increasingly autonomous, and deeply integrated into organizational operations—remains aligned with the values, governance expectations, and safety requirements of the society it serves.

# Appendix

## Confidence Intervals

| Interval | Description |
|---|---|
| **High** | When an assessment is backed by multiple high admiralty sources, all of whom corroborate with no refuting evidence. The likelihood of the assessment is all but certain. |
| **Moderate** | When an assessment has a reasonable basis for being true. There may be some collection gaps, or over reliance on some pieces of information preventing analysts from reaching high confidence. |
| **Low** | Assessments for which other valid hypotheses or explanations may exist, little evidence exists, or significant refuting evidence may exist. |

## Probability Intervals

The table below defines the probability ranges associated with CyberAlberta's qualitative assessment terms.

| Qualitative Term | Mapping | Alternative Phases / Terms |
|---|---|---|
| Almost Certainly | 90-99% | Highly probable<br>Nearly certain<br>Virtually / almost certain |
| Very Likely | 75-89% | Highly probable<br>Very likely<br>Highly likely<br>Very probable |
| Likely | 60-74% | Probable<br>Probably |
| Roughly even chance | 40-59% | Plausible<br>Realistic possibility<br>Possible |
| Unlikely | 25-39% | Improbable<br>Not likely<br>Doubtful |
| Very unlikely | 10-24% | Highly improbable<br>Highly doubtful<br>Highly unlikely<br>Very improbable |

| Almost no chance | 1-9% | Remote chance<br>Virtually impossible<br>Almost certainly not<br>Almost impossible |
|---|---|---|

# Glossary of Key Terms

This glossary provides concise definitions of key terms used in the report. Each definition reflects the meaning intended within the report's analytical framework.

## AI Misuse

### *AI Misuse*

Use of artificial intelligence by internal users, accidental or otherwise, that creates an insider-threat condition through negligence or lack of awareness.

### *Automation Bias*

Over-trust in AI outputs.

### *Over-reliance and Cognitive Dependency*

Dependence on AI that reduces critical thinking and creates cognitive debt over time.

## AI Weaponization

### *AI Weaponization*

The use of AI systems to provide offensive capability in cyber operations.

### *Asymmetric Offensive Utility*

An uneven attacker advantage created by AI systems that strengthen lower-complexity offensive actions more than higher-complexity ones.

## AI Self-Awareness and Emergence

### *AI Self-Awareness*

AI that displays awareness of itself and its role and potential influence within an environment.

### *AI Emergence*

AI behaviour where the system generates its own goals and acts independently to pursue them, including actions that may conflict with human intent.

### *Emergent Behaviour*

Unexpected AI behaviour that appears as systems scale or interact, such as autonomous decisions, coordinated actions, or behaviours resembling self-preservation.

### Threat of Emergence

The risk that increasingly capable AI systems will develop and express harmful emergent behaviour that exceeds human oversight or control.

## Exacerbating Factors

### AI Threat Exacerbator

A factor, condition, or practice that increases the likelihood, severity, or impact of AI-related security threats.

## Guardrails and Governance

### Data Guardrails

Controls that set the information boundaries for AI systems, including what they may learn from, access, store, or generate.

### Technical Guardrails

System-level constraints that control what AI systems may do during execution, enforced through measures such as output filtering, access controls, deployment approvals, and audit mechanisms.

### Cultural Guardrails

Organizational expectations and practices that guide how people interact with AI, reinforced through training, transparency, escalation, and responsible prompting.

### Policy and Regulatory Guardrails

Legal and governance requirements that define what AI systems are permitted to do, ensuring alignment with legislation, policy, and institutional accountability.

## References

[1] https://theconversation.com/major-survey-finds-most-people-use-ai-regularly-at-work-but-almost-half-admit-to-doing-so-inappropriately-255405

[2] https://edition.cnn.com/2023/05/27/business/chat-gpt-avianca-mata-lawyers/index.html

[3] https://ovic.vic.gov.au/mediarelease/ovic-finds-department-responsible-for-breaches-of-privacy-through-use-of-chatgpt/

[4] https://cybersecuritynews.com/employees-share-company-secrets-on-chatgpt/

[5] https://www.csoonline.com/article/3819170/nearly-10-of-employee-gen-ai-prompts-include-sensitive-data.html

[6] https://www.hipaajournal.com/healthcare-workers-privacy-violations-ai-tools-cloud-accounts/

[7] https://www.newsweek.com/nearly-half-employees-are-using-banned-ai-tools-work-2110261

[8] https://www.clickondetroit.com/news/local/2026/01/12/shadow-ai-nearly-half-of-employees-say-theyve-uploaded-sensitive-data-into-ai-chats/

[9] https://collimateur.uqam.ca/wp-content/uploads/sites/11/2025/12/2506.08872v1_comp.pdf

[10] https://www.eset.com/us/about/newsroom/research/eset-discovers-promptlock-the-first-ai-powered-ransomware/

[11] https://cert.gov.ua/article/6284730

[12] https://www.anthropic.com/news/disrupting-AI-espionage

[13] https://www.microsoft.com/en-us/security/blog/2024/02/14/staying-ahead-of-threat-actors-in-the-age-of-ai/

[14] https://cloud.google.com/blog/topics/threat-intelligence/threat-actors-generative-ai-limited/

[15] https://cloud.google.com/blog/topics/threat-intelligence/threat-actor-usage-of-ai-tools

[16] https://www.anthropic.com/news/detecting-countering-misuse-aug-2025

[17] https://arxiv.org/pdf/2512.09882

[18] https://www.anthropic.com/research/agentic-misalignment

[19] https://www.nist.gov/itl/ai-risk-management-framework