# A Framework for Safeguarding Generative Al and Agentic Al Models with Foundational Security Principles



# The Challenge:

The increasing adoption of both Generative and Agentic Artificial Intelligence (AI) models presents significant opportunities for innovation and efficiency. However, these powerful models also introduce new security vulnerabilities. Organizations face the critical challenge of protecting these valuable assets from evolving threats, including unauthorized access, data breaches, model manipulation, and intellectual property theft. Traditional security approaches have proven to be insufficient to address the unique risks associated with protecting AI application models.

## The Solution: Foundational Security for Al

To effectively safeguard generative and agentic AI models, organizations must implement a robust and multi-layered security strategy built upon established and proven principles. This solution brief highlights the critical role of encryption, Zero Trust Architecture (ZTA), stringent access controls as foundational elements for securing these advanced AI systems, as emphasized by leading authorities in cybersecurity.

### **Key Components of the Solution (Supported by Authoritative Guidance):**

- Robust Encryption: Implementing strong encryption methods, such as two-way encryption (mTLS) is paramount for protecting generative and agentic AI models and their associated data, both at rest and in transit. This ensures that sensitive information remains confidential and inaccessible to unauthorized parties, even in the event of a breach. This principle is consistently recommended in cybersecurity best practices and highlighted in guidance from organizations like the National Institute of Standards and Technology (NIST) and the Cloud Security Alliance (CSA).
- Zero Trust Architecture: Adopting a Zero Trust security model is crucial for mitigating risks within the AI environment. This approach operates on the principle of "never trust, always verify," requiring strict end-source root-signed certificate authentication and authorization for every user, device, and application attempting to access AI models and related resources. This is a key recommendation in NIST's cybersecurity frameworks and is increasingly recognized as a vital strategy for modern security, including AI, as noted by CSA.
- Stringent Access Controls: Implementing granular and consistently enforced access control policies is essential for limiting who can interact with AI models and their underlying data. This includes certificate-based access control (CBAC), and continuous monitoring of access patterns to detect and prevent unauthorized penetration or modification. The Open Web Application Security Project (OWASP) AI Security Project and guidelines from the European Union Agency for Cybersecurity (ENISA) underscore the importance of robust access management for AI systems.

# **Benefits of Implementing Foundational Security:**

By prioritizing encryption, Zero Trust Architecture, and stringent access controls, organizations can:

Reduce the Risk of Data Breaches: Protect sensitive training data and model outputs from unauthorized access and
exfiltration (as highlighted by OWASP and industry threat reports).

- **Prevent Model Manipulation:** Ensure the integrity and trustworthiness of AI models by controlling who can access and modify them (a key concern addressed by NIST).
- **Enforce Data Integrity for MLM:** Prevent the manipulation of machine-learning model (MLM) data by controlling the data source to avoid altering machine-learning algorithms negatively.
- **Protect Intellectual Property:** Safeguard proprietary AI algorithms and model weights from theft or reverse engineering (a focus in discussions around AI security best practices).
- **Enhance Regulatory Compliance:** Meet increasingly stringent data protection and security requirements related to AI (relevant to regulations like the EU AI Act).
- **Build Trust and Confidence:** Demonstrate a commitment to security, fostering greater trust in the organization's Alpowered products and services.

### Conclusion:

Securing generative and agentic AI models is a critical imperative for organizations seeking to leverage their transformative potential responsibly. Implementing foundational security principles such as encryption, Zero Trust Architecture, and stringent access controls, as advocated by leading cybersecurity authorities like NIST, OWASP, and CSA, provides a robust framework for mitigating key risks and ensuring the confidentiality, integrity, and availability of these valuable AI assets. By prioritizing these fundamental measures, organizations can confidently embrace the power of AI while safeguarding their operations and the trust of their stakeholders.

# **Trustus Advance Intelligence Security**

Trustus AI Security is an advanced unified digital identity cybersecurity and privacy platform delivering innovative internet communications. We partner with businesses of all sizes to provide automated trusted solutions for securing digital identities, and connecting devices, applications and ecosystems, enabling people and things to communicate with each other safely. Our team has delivered some of the most complex web portal and infrastructure solutions to some of the best-in-class international companies. Managing the complexities and risks associated with internet communications is our strength. The Trustus team looks forward to work with you.

Contact our team of experts today to discuss your specific needs at **Sales@TrustusSecurity.com**. Get in touch.

https://trustussecurity.com/

Legal Disclaimer: This presentation is issued by Trustus Technologies Group (Trustus). All software technologies architected, designed, developed, customized, and deployed; current and future, including to be and not yet deployed; all content as herein described are the exclusive proprietary intellectual property (IP) of Trustus. Any content herein, views or opinions express or implied, are our own and are offered without any fault or recourse by the readers. Readers are cautioned they cannot use the information contained herein as their own, nor engineer or reverse engineer the IP without the express written permission of Trustus.

