Proprietary + Confidential

Google Cloud

Protecting the future of Al

How Google Cloud helps mitigating LLM risks and promote a responsible adoption

April 7th 2025





Responsible Al is



Our work is guided by the Al Principles

Al should:

2

3

5

6

- Be socially beneficial
- Avoid creating or reinforcing unfair bias
- Be built and tested for safety
 - Be accountable to people
 - Incorporate privacy design principles
 - Uphold high standards of scientific excellence
 - Be made available for uses that accord with these principles

Applications we will not pursue:

Likely to cause overall harm

Technologies primarily intended to cause injury

Surveillance violating internationally accepted norms

3

Purpose contravenes international law and human rights

Al Governance Committee



Three ways we build security into the Google Cloud



Google's Secure Al Framework (SAIF)

Al is advancing rapidly, and it's important that effective risk management strategies evolve along with it



Expand strong security foundations to the Al ecosystem



Extend detection and response to bring Al into an organization's threat universe



Automate defenses to keep pace with existing and new threats



Harmonize platform level controls to ensure consistent security across the organization



Adapt controls to adjust mitigations and create faster feedback loops for Al deployment



Contextualize Al system risks in surrounding business processes

SAIF Risk Map

https://saif.google

Application

- 1. Denial of ML Service
- 2. Insecure Integrated Component
- 3. Model Reverse Engineering
- 4. Rogue actions

Model

- 5. Sensitive Data Disclosure
- 6. Inferred Sensitive Information
- 7. Prompt Injection
- 8. Model Evasion
- 9. Insecure Model Output

Infrastructure

- 10. Excessive Data Handling
- 11. Model Source Tampering
- 12. Model Exfiltration
- 13. Model Deployment Tampering

Data

- 14. Data Poisoning
- 15. Unauthorized Training Data



Model Tuning

- Parameter Efficient Fine Tuning (PEFT)
- Customer specific adapter weights
- Foundational model remains frozen during inference



- Input data is secured at every step
- Adapter weights are stored securely
- Customer can delete adapter weights at any time
- Customer Data will not be logged to train foundation models by default



Google Cloud Architecture





Takeaways

- Establish Al Governance, and build Responsible Al capabilities
- Explore AI development through a security lens
- Al security has to be addressed as a **company-wide** challenge
- Leverage an **Google SAIF framework** to make sure you have a comprehensive view of all the AI risks
- Implement technical measures to protect Al applications, using cloud-native controls, like Model Armor



Google Cloud

Proprietary + Confidential

Thank You!

in linkedin.com/in/dluzi



