

# Expect the unexpected

Incident response strategies for GenAI

Jolanda Kumakaw | Manager, Strategic Command, Google Trust & Safety



# Topics

- 01 What are GenAI incidents?
- 02 Emerging risk areas
- 03 Incident response
- 04 Examples
- 05 Q&A



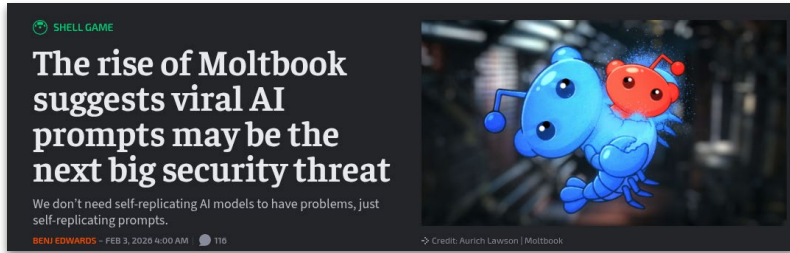
**“ P0 incident incoming!”**



# Google's 'Woke' Image Generator Shows the Limitations of AI

Google has hit pause on Gemini's ability to generate images of people after a far-right backlash to its historical depictions.

Wired, Feb 22, 2024

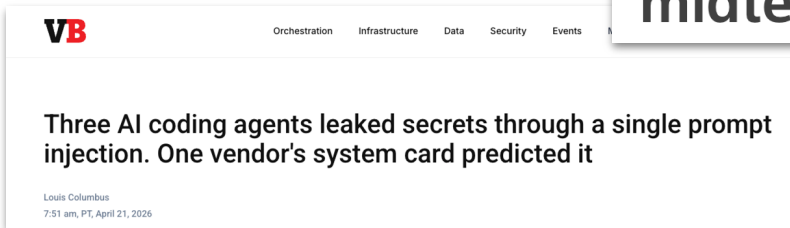


Grok 'nudify' scandal: Data Protection Commission to investigate X over its AI app

AI chatbots can sway voters with remarkable ease – is it time to worry?



AI deepfakes blur reality in 2026 US midterm campaigns



Anthropic's Mythos sends US banks rushing to plug cyber holes

Reuters, May 12, 2026

# What are “GenAI incidents”?



## Cybersecurity incidents

Focus on *confidentiality, integrity, and availability* of the system. Can affect any product.

### EXAMPLES:

- Unauthorized access
- Malware
- DDoS attacks



## AI incidents

Focus on the *output and logic* of the model. Occur *on* or *via* AI products/models.

### EXAMPLES:

- **On AI:** Hallucinations, safety risk outputs
- **Via AI:** Generated misinformation (e.g., NanoBanana)



## AI-Connected Cyber incidents

Cybersecurity incident that has an AI connection where AI acts as a **multiplier**, scaling traditional threats and increasing safety risks.

### EXAMPLES:

- AI-enabled automated phishing
- Adversarial prompt injection
- Data exfiltration via email tools

# Emerging risk areas

mov  
ie\_fil  
ter

## Better GenMedia

GenMedia producing better synthetic media.

- How does this impact upcoming elections?
- Nudification of people without their consent (NCII)

pers  
on

## AI Personalization

Users are able to have a personalized experience based on their data and interactions with AI.

- Users may think their personal info is compromised.
- Where is the “creepy line” for users?

sma  
rt\_t  
oy

## Agentic Era

Entering the agentic era - autonomous tasks resulting in less human in the loop (HITL).

- Agent makes incorrect inference and executes the wrong action.
- Prompt injections bypass security or results in user harm

# When parachuting into chaos during a GenAI incident, remember to:



## Connect

**Comprehend** the issue and risks.  
Identify **cross product** impact and trends.  
Identify **gaps** and **connect the dots**.  
Be the **connective tissue** between distinct response efforts.



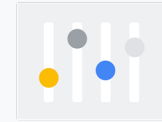
## Coordinate

**Organize** tactical response across distinct teams.  
**Assemble** response teams.  
Lead response **strategy**.  
Be the **single-threaded** owner.



## Communicate

- Central hub** of all critical incident information.
- Owner of primary incident communication channels.  
**Owner** of **source of truth** incident real-time doc.
  - Notify Execs and key stakeholders.



## Control

- Steer** the response effort. Jump in.  
Set **priorities** based on **contextualized** understanding of risks.  
**Ensure alignment** across teams.  
Set up necessary **adhoc workflow/streams**
- **Accelerate** where needed

# Examples

**Nudification  
using GenAI**



Personalized  
user  
experience in  
GenAI  
products

# Q & A

Get in contact via LinkedIn: [/jolanda-kumakaw](#)