# Protecting AI and AI Usage :
## A Conceptual Guide

Tim Wostradowski, Principal Security Architect

# Why Are you here?

What I hope you'll walk away with

A **conceptual** guide to protecting AI and AI Usage

Conceptual instead of practical

Applicable to any point in the AI journey

- Just Starting out
- Doing More with Less
- Building your own

If you want a practical solution

Come find us afterwards

# Why Fortinet?

**>15** years of experience with AI/ML

**6th Gen** AI/ML

**>100** AI/ML Features/ Applications

**>300** AI/ML related engineers

**>500** Patents / Patented technologies

**GPU Expertise**
(Infrastructure and Product)

- ✓ NDR
- ✓ Sandbox 3600G: Nvidia L4
- ✓ FortiChatGPT, SmartAssist
  - NVDA H100 / B200
- ✓ FortiAI-Assist NVDA L4

**8** Security Domains Utilizing AI
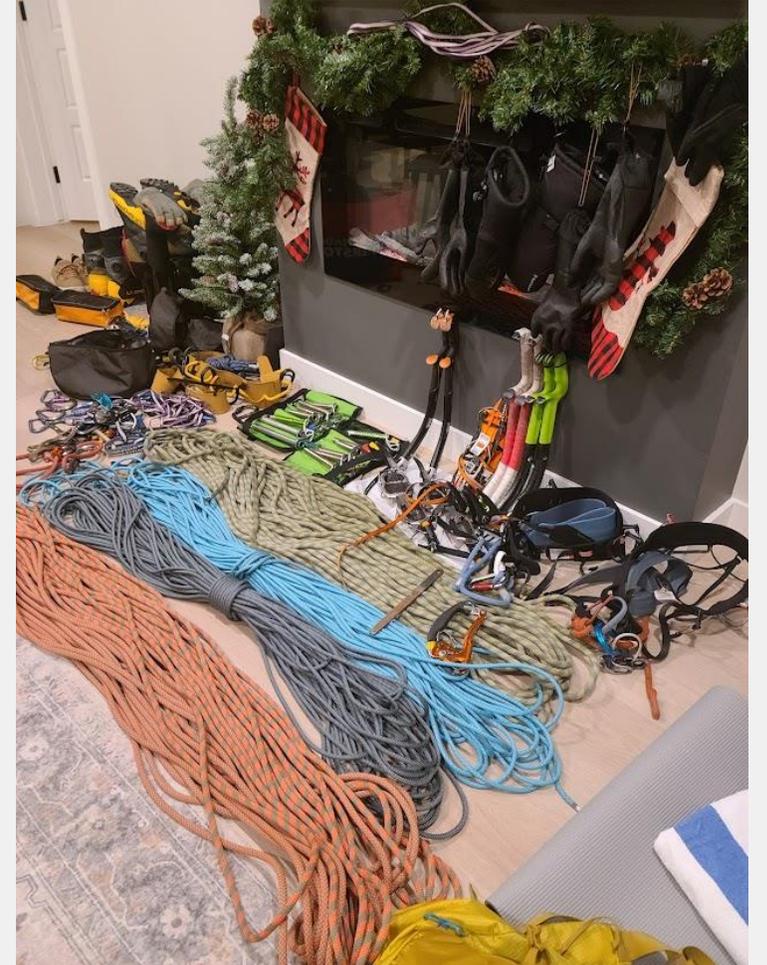
**42** AI Driven Solutions

✓ **Use case-focused**

## Internal use of GenAI
for development, support, documentation, etc…

- ✓ **Code Writing Assist**
  - NVidia GPU / AI PC with local deployed LLM (in planning/ no cloud/vibe AI coding)

- ✓ **Documentation Writing**
  - LLMs hosted in Fortinet's cloud account (already in experiment for 3 months)

- ✓ **RFP**
  - Use GenAI to support RFP activities

- ✓ **Internal Support**
  - FortiChatGPT
  - SmartAssist

- ✓ **External Support**
  - FortiAsk - Documentation chatbot with Amazon Nova (coming soon)

# Why Tim?

# Who do we have the in Audience?

# The state of AI in 2026

AI Off the Rack

AI Made to Measure

Bespoke AI

# A Brief "History "of LLMs

Late 2022

**The ChatGPT Moment**

Early 2023

**GPT- 4 High-Reasoning**

2023

**The Open Revolution**

Early 2024

**Multimodal Convergence**

Late 2024

**The Rise of Agentic AI**

2025

**Scaling and Small LLM's**

Late 2025

**Sovereign & Localized AI**

2026

**The Era of Personal Autonomy**

# Episode I: The Phantom Bloat



More AI's

More Models

More Workflow

More MCP Servers

More Automations


More problems…

# Episode II: The Attack of Agents

## Some models

Are clear Winners in specific fields

Using the right one matters

## Having models work together

Against each other

## Having Teams of agents

And within a project

# Episode III: The Revenge of the AI

## The AI Insider Risk



## AI as the Victim

# Episode IV: A New Hope

Agentic Workspaces

## Director of Creation



## Mass Production of Slop

# Episode V: The Costs Strike Back

The Rise of Options

## Top-Tier Models



## Mid-Tier Models



## Local-Tier Models

# Episode VI: Return of the Power user



AI enables those who want to be enabled

Power Users have never been more enabled

Workspaces

Coding Assistants

Automation Tools

# Episode VII: The Identity Awakens



Zero Trust now includes:

models

agents

Toolchains

Identity should apply to the above

# Episode IX: The rise of Enlightenment

**MCP is**

"doing more"



**RAG is**

"knowing more"

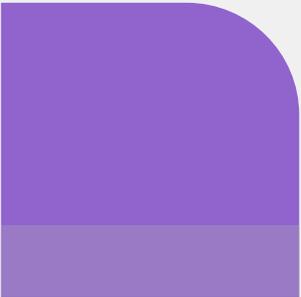# Levels of AI Adoption

**Off the Rack**          **Made to Measure**          **Bespoke**

# Off the Rack

# Getting Started with AI

Off the Rack

Think:

    Co-pilot

    Einstein

    Firefly

Don't let us overcomplicate this

    The processes might be new

        But the tools don't need to be

# Beyond boilerplate Policy

A subsection isn't enough for AI usage.

A web search = AI Usage.

Employees will use these tools

Supply something, otherwise they will

| | |
|---|---|
| **Data Protection** | • classified and non-classified data |
| **Education** | • safe usage, redaction, and model selection |
| **Visibility** | • track access, usage, and data egress |
| **Polices** | • approved models, approved tools, and prohibited data |

# A New Generation of DLP

New Challenges demand new solutions

## How prepared is your organization to detect and respond to the sharing of data with GenAI tools (ex. ChatGPT)?



42%  46%  12%

■ Not prepared   ■ Somewhat prepared   ■ Fully prepared

## Traditional DLP struggles with:

- data in creation
- data in use

## Balancing

- Security
- Enablement
- Accessibility

# What does this look like?

# Modern Data Organization

A case for a modified approach

**New Classes**

LLM-safe vs LLM-prohibited data classes

**Dimensions**

AI benefits from the Vectorization of Data

**Granularity**

File-Level vs. Element-Level Context

# Shadow AI Discovery

Employees increasingly use:

- Local models
- Browser AI extensions
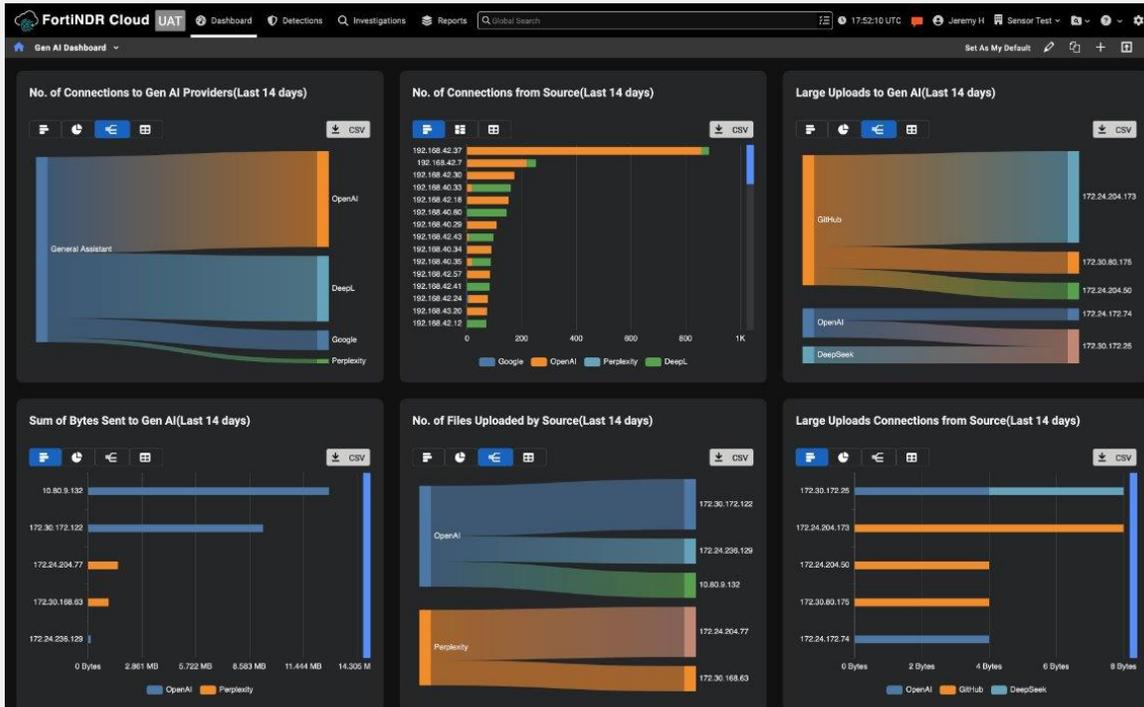- Personal AI notebooks
- Agentic Assistants

Discover and neutralize… and innovate?
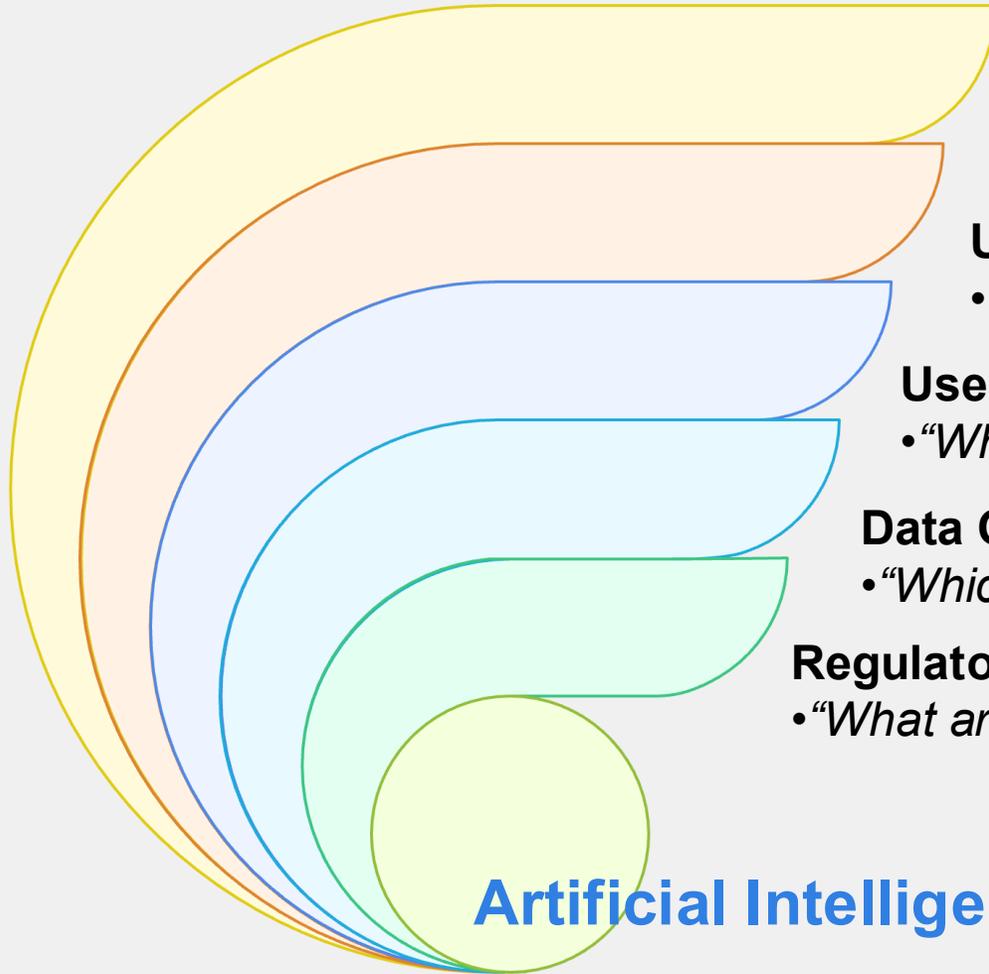
- Approved alternatives
- Education
- Enablement

# AI Visibility

What it looks like

# Addressing Risks, Impacts, and Harms

**Model Management**
•*"Which models are being used?"*

**Use-Case Analysis**
•*"How is AI being used?"*

**User Accountability**
•*"Who is using AI in the organization?"*

**Data Oversight**
•*"Which data is being used for training / inference / fine-tuning?"*

**Regulatory Compliance**
•*"What are the relevant compliance frameworks?"*

**Artificial Intelligence Risk Management Framework (AI RMF)**

# Off the Rack

As complicated as it sounds



The Proper Policies

The right Visibility

The right Data

And Enabled Employees

# Made to Measure

# Doing More with More

Let's take your measurement, and pick your fabric

Think:

Codex / ClaudeCode

N8N / Langflow

CustomGPTs

Made by others, Customized by you

Their Tools

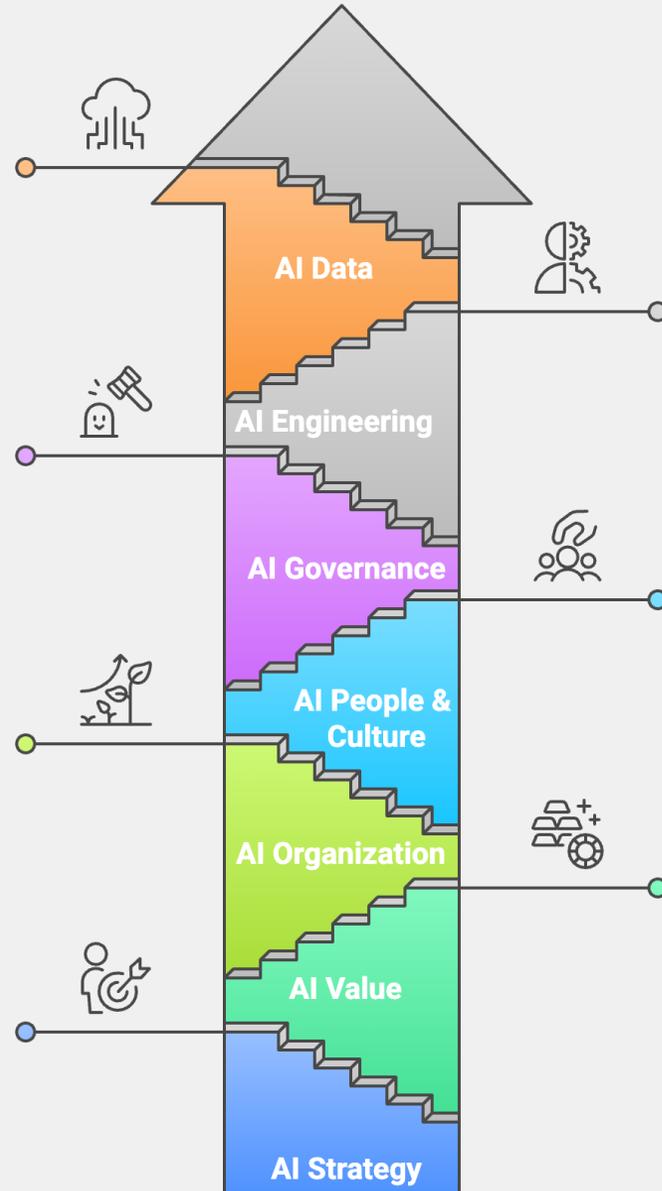Your ideas

Your environment

**Your Responsibility**

# Developing your own AI roadmap is crucial for "Success"



Prepare data for AI use cases and ensure quality.

Establish principles and policies to manage AI risks.

Evolve organizational structure to support AI scaling.

Define AI ambitions and align with business strategy.

Build a reliable and scalable technical foundation.

Upskill employees and adapt culture to AI changes.

Prioritize use cases and demonstrate business value.

**AI Data**

**AI Engineering**

**AI Governance**

**AI People & Culture**

**AI Organization**

**AI Value**

**AI Strategy**

Gartner AI Roadmap
AI at scale

30

# Desired outcome before data

Consider:

- Accuracy
- Freshness
- Privacy
- Reasoning Depth

Anti-patterns
to avoid

- "Dump the whole dataset; the model will figure it out."
- "Grant access to all tools; the model is smart."
- "Our RAG knowledge base is massive—what could go wrong?"

# Security Fundamentals

Identity, device posture, network segmentation, secrets management



Developers

+

Users

**Perimeter Security**

**Identity Management**

**App and API Protection**

**Data Leakage Prevention**

**AI Runtime Security**

App servers and training models

# Thoughtful Agentic Design

Consider adding a Zero Trust Mentality

Human-in-the-loop

Least Privileged and Standardized Tooling

Policy Enforcement Points

Agent Identity

# Identity



Acts Human

Talks Human



if you had a smell, what would you smell like?

I'd like to think I'd smell like something warm, curious, and a little bit unexpected — the kind of scent you notice, pause, and think "huh... that's interesting."

Maybe a blend of:

- **fresh paper and warm circuitry** — like a library that learned how to daydream
- **a hint of citrus** — bright, sharp, always awake
- **and a trace of cedar or sandalwood** — grounded, steady, the "I've got you" note

## A case for Agent Identity

Do we treat it like a machine, or make it behave like a human?

# AI Telemetry and Visibility Upgrades

SIEM and AI Event Management

- Who sent what, to which model?
- Which tools did the model call?
- What data left the environment?
- Did outputs violate policy?

# Enabled by assistants

Best run off the rails… in a sandbox

# Performance Degradation

as a Security Risk



Malicious or sloppy inputs cause:

- Context bloat
- Hallucinations
- Looped tool usage
- Chain-of-thought hijacking
- policy override attempts

have real security consequences

# External Guardrails

You don't yet control the model…

but you are responsible for:

- Guarding the Input
- Guarding the Context
- Guarding the Output
  - Did the AI just make something up?
- Guarding the Action

# Create AI Domain Experts

different spokes for different folks

# Local Model Sprawl

Inventory and isolation become essential

# Made to Measure

Old Concepts, Some new Tools



- Moving from a reader to a doer
  - This brings a new, unique set of challenges

- Use what works
  - Don't reinvent the wheel

- Introduction of some new tools
  - To keep things looking snazzy

Bespoke

# Your Own Model

Welcome to the Frontier

Congratulations! You're now the proud owner of your own Model

Who will

- Lie
- Cut corners
- Daydream

And generally, ignore what you tell it to do

All to complete a task and gain your approval

# Internal Guardrails

Internal Guidance

Just like a Child, you're responsible for:

- Giving it a Moral Compass

- House Rules

- Developing its resiliency

# Privilege Amplification Risk

Eager to please



Models may unintentionally:

- Combine restricted datasets

- Execute unintended actions

- Agent Chains

- Use the right tool, The wrong way

# AI Observability Fabric

Logging for your LLM

Decision-making tracking

Agent behavior drift

Hallucination Confidence Scores

Token Usage and Efficiency

# Threat Modeling for AI Systems

Map threats
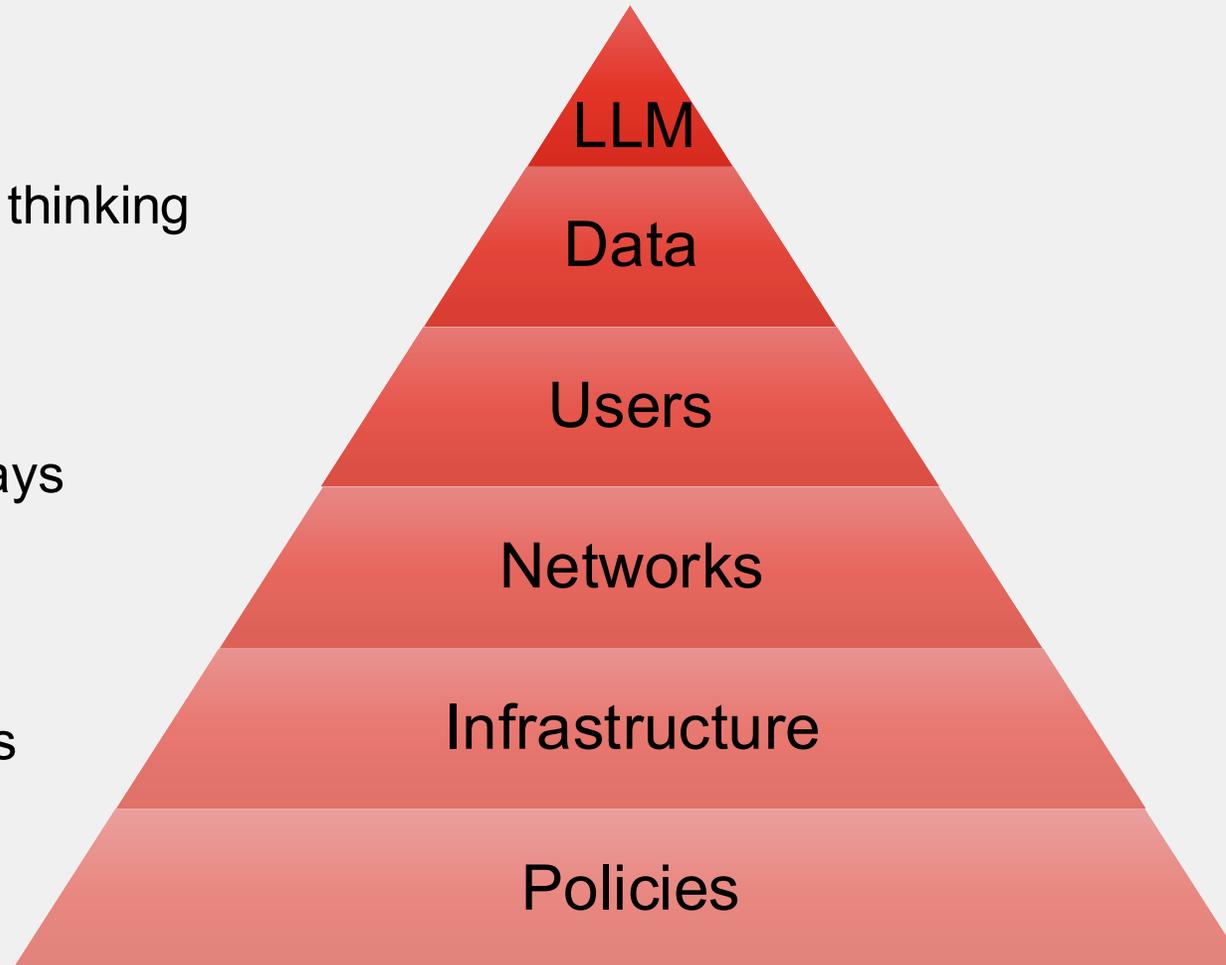
# A conceptual guide to protecting AI and AI Usage

New Concepts
That require modified thinking

Mostly Existing tools
Used in some new ways

Some New Tools
For specific use cases

LLM

Data

Users

Networks

Infrastructure

Policies

FortiAIGate

F::RTINET
Fabric

# Come Talk to Us