# $whoami

Currently: Chief AI Safety Officer/Co-Founder at **KNOSTIC**

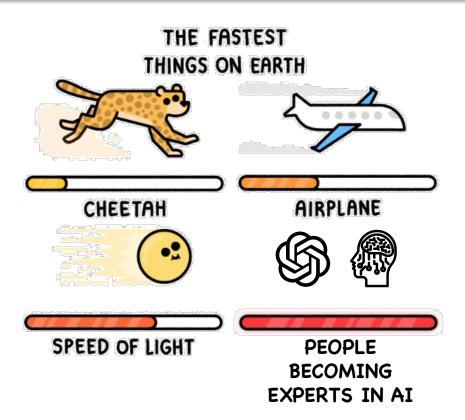Formerly: Chief Security Scientist at **BANK OF AMERICA**

Other Random Recognitions:
- Lifetime Achievement Award
- Cybersecurity Hall of Fame
- Most Influential Person in Security
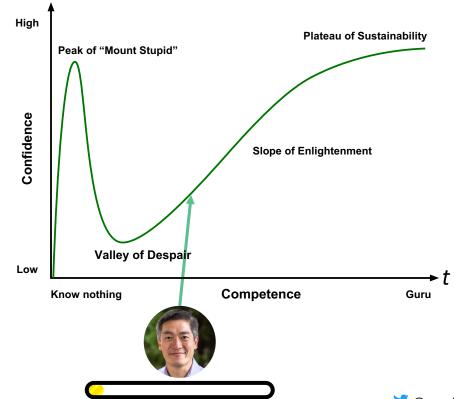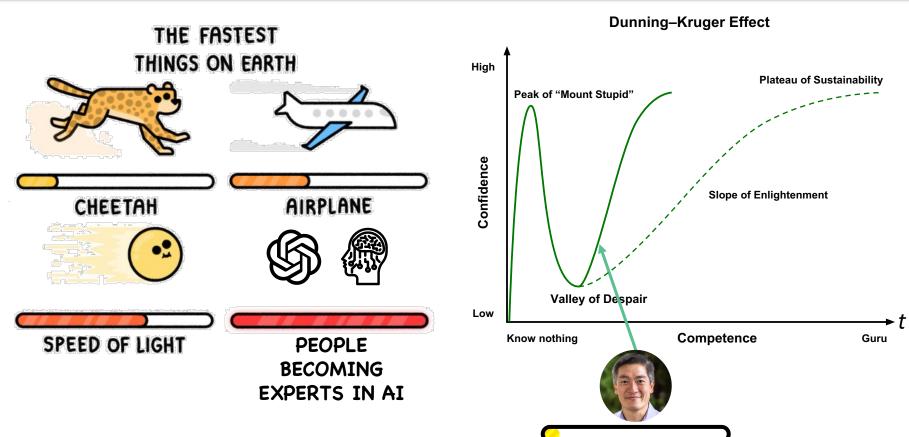- Cyberscoop Cybersecurity Visionary
- 20+ Patents

THE FASTEST
THINGS ON EARTH

CHEETAH

AIRPLANE

SPEED OF LIGHT

PEOPLE
BECOMING
EXPERTS IN AI

**Dunning–Kruger Effect**



High

Confidence

Low

Peak of "Mount Stupid"

Plateau of Sustainability

Slope of Enlightenment

Valley of Despair

Know nothing

Competence

Guru

$t$

KNOSTIC

@sounilyu

# How do you accelerate competence?

THE FASTEST THINGS ON EARTH

CHEETAH

AIRPLANE

SPEED OF LIGHT

PEOPLE BECOMING EXPERTS IN AI

Dunning–Kruger Effect

Confidence

High

Peak of "Mount Stupid"

Plateau of Sustainability

Slope of Enlightenment

Valley of Despair

Low

Know nothing

Competence

Guru

t

KNOSTIC

@sounilyu

# Accelerating Competence w/ Mental Models

AI for Good
Global Summit
11 July 2025
Geneva, Switzerland

**Cognitive Biases**
Our mind has well documented biases that lead to us making irrational thoughts and decisions

**Behaviour Change / Persuasion**
How do you influence yourself and others?

**Models & Data Analysis**
Find and display insights

**Leadership**
Bringing people along to your vision of a bigger and brighter future

**Communication**
Getting your message clearly across to your audience

**Project management**
Know Problem; Know Solution
Know Problem; Unknown solution
Unknown Problem; Unknown Solution
The project management style you pick is critical to the successful outcome

**Goal Setting**
Tools to help you set the right goal

**Risk Management**
Limiting the downside and maximise the upside

THE 10X ENGINEER – Mental Models
Tom Connor 2018
https://medium.com/10x-curiosity

**Finance**
Skills to manage and understand money
- Balance sheet
- Cash Flow Quadrant
- Compound interest
- Externality
- NPV
- Opportunity cost
- Pay yourself first
- Payback period
- Sunk Cost
- Supply and Demand
- Index Funds
- EBITA

**Understanding yourself and others**
We are all different

**Problem Solving**
Look for problems to drive change. There is abundance in the world and we can both win

**Scientific Method**
If P is Low the null must go!

**Time Management / Execution**
How do you deliver your priorities relentlessly?

**Prioritisation (80/20)**
What are the few key task you can do to delivery the majority of the value?

KNOSTIC

@souniyu

# Cyber Defense Matrix in 3D!



AI for Good
Global Summit
11 July 2025
Geneva, Switzerland

KNOSTIC

@sounilyu

# Cyber Defense Matrix + Data?

KNOSTIC

@sounilyu
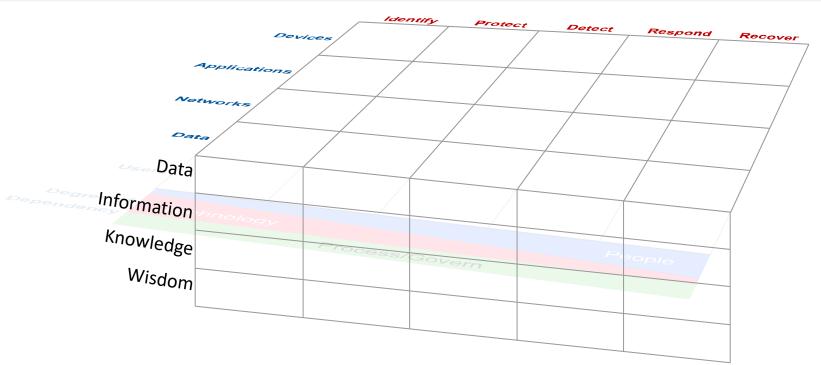
# Cyber Defense Matrix + DIKW
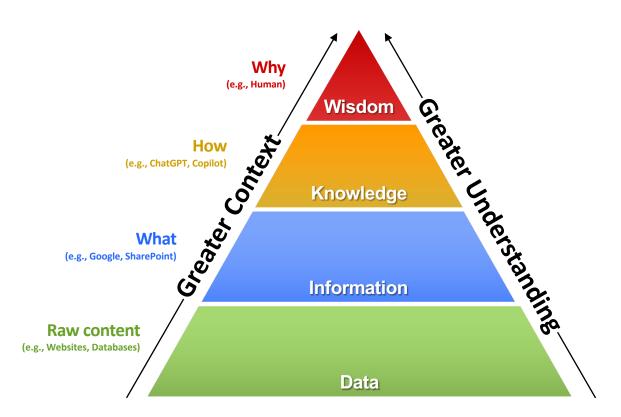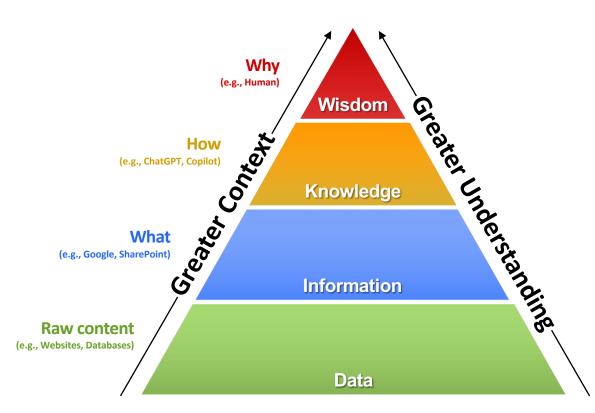
# The DIKW Pyramid

ChatGPT has unlocked the knowledge economy that will power the AI era
- Knowledge Engineering
- Knowledge Lakes
- Knowledge Pipelines
- Knowledge Quality
- Knowledge Security
- Knowledge Governance
- Knowledge Search
- Knowledge Sharing
- Knowledge-based Decision Support
- Knowledge Classification
- Knowledge Provenance
- Knowledge Management
- Knowledge Privacy

# Access Controls

Why
(e.g., Human)

How
(e.g., ChatGPT, Copilot)

What
(e.g., Google, SharePoint)

Raw content
(e.g., Websites, Databases)

Wisdom

Knowledge

Information

Data

Greater Context

Greater Understanding

???

Information-centric controls
(e.g., classification)

Data-centric controls
(e.g., file level permissions)

KNOSTIC

@sounilyu

# What Happens When Applying Old Controls



AI for Good
Global Summit
11 July 2025
Geneva, Switzerland

**Why** (e.g., Human)

**How** (e.g., ChatGPT, Copilot)

**What** (e.g., Google, SharePoint)

**Raw content** (e.g., Websites, Databases)

Greater Context

Greater Understanding

Wisdom

Knowledge

Information

Data

???

Information-centric controls (e.g., classification)

Data-centric controls (e.g., file level permissions)

KNOSTIC

@sounilyu

# What Happens When Applying Old Controls

# Inference Is Still a Problem

**Why**
(e.g., Human)

**How**
(e.g., ChatGPT, Copilot)

**What**
(e.g., Google, SharePoint)

**Raw content**
(e.g., Websites, Databases)

Greater Context

Greater Understanding

Layoffs

Wisdom

Knowledge

Information

Future Building Space

Future Equipment Purchases

Data

Even if your data permissions are **perfect**, you still have inference problems

Information-centric controls
(e.g., classification)

Data-centric controls
(e.g., file level permissions)

KNOSTIC

@sounilyu

# Knowledge Centric Capabilities → Knowledge Centric Controls

**Why** (e.g., Human)
**How** (e.g., ChatGPT, Copilot)
**What** (e.g., Google, SharePoint)
**Raw content** (e.g., Websites, Databases)

Greater Context

Greater Understanding

Wisdom

Knowledge

Information

Data

*Knowledge*-centric controls (e.g., need-to-know)

Information-centric controls (e.g., classification)

Data-centric controls (e.g., file level permissions)

KNOSTIC

@sounilyu

KNOSTIC

@sounilyu

# Takeaways From the DIKW Pyramid

**Why** (e.g., Human)
**How** (e.g., ChatGPT, Copilot)
**What** (e.g., Google, SharePoint)
**Raw content** (e.g., Websites, Databases)

Greater Context

Greater Understanding

Wisdom
Knowledge
Information
Data

1. Recognize when a new layer of abstraction has appeared

2. The macro patterns may be similar, but the specific standards will not apply across layers

3. Align problems and solutions at the same layer

4. Lower-level controls can negatively impact higher-level capabilities

5. Higher-level controls can make lower-level problems less relevant

KNOSTIC

@sounilyu

# How YOLO Do You Want to Go?

**Automated Patching, Threat Intel Blocklists**

**LLMs + Tools, SOAR**

**LLMs + Tools + Agentic AI**

| Sensing | Sensing | Sensing |
|---|---|---|
| Sense Making | Sense Making | Sense Making |
| Decision Making | Decision Making | Decision Making |
| Acting | Acting | Acting |

**Useful Controls**
- Reliable sensors (identity)
- Brakes and reverse gear
- Narrowly scoped actions

**+**

**Useful Controls**
- Regression testing
- Validated assumptions
- Documented processes

**+**

**Useful Controls**
- Pre-established thresholds
- Pre-determined authorities
- Clear lines of accountability
- Imagine 100 interns, WCPGW?

KNOSTIC

@sounilyu

# Summary

**Knowledge-Centric Capabilities Require Knowledge-Centric Controls**



**Many Opportunities for Standards Even Before Agentic AI**



**There's a Right and Wrong Time for Standards**



**Span of Control for Agents Depends on Complexity and Value**