# Assessing "AI-Enabled" Tools

## Separating the Wheat from the Chaff

**Mike Hadjimichael, Ph.D.**

**Cyber AI Domain Capability Area Lead, National Cybersecurity FFRDC**

**The MITRE Corporation**

**mikeh@mitre.org**

**MITRE** | SOLVING PROBLEMS FOR A SAFER WORLD™

# There is Still Much to Learn With AI



Photo: Cornell University

Image from: https://spectrum.ieee.org/cars-that-think/transportation/sensors/slight-street-sign-modifications-can-fool-machine-learning-algorithms
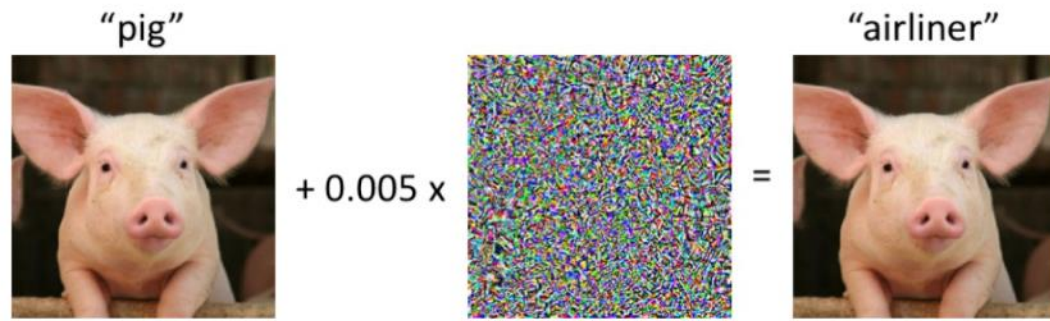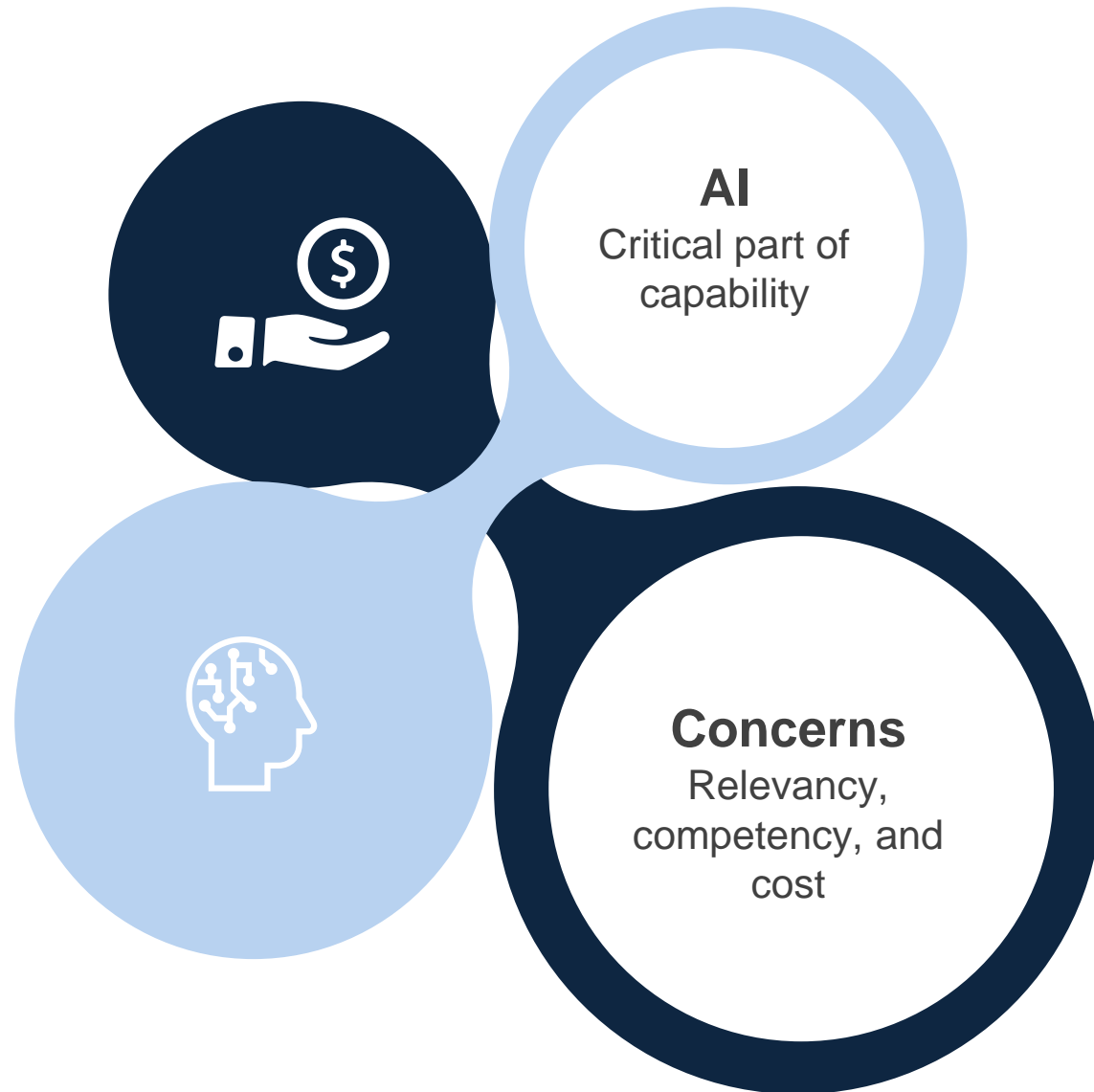


"pig"  + 0.005 x  =  "airliner"

Image from https://gradientscience.org/intro_adversarial/

MITRE

# Challenge

**AI**
Critical part of capability

**Concerns**
Relevancy, competency, and cost

**Uncertainty**
What questions do we ask to pair functionality of tool with what AI does?

**MITRE**

# Addressing the Challenge



https://www.nist.gov/cyberframework

https://shield.mitre.org

**Frameworks**

CORE

| Identify-P |
| Govern-P |
| Control-P |
| Communicate-P |
| Protect-P |

https://www.nist.gov/privacy-framework

**ATT&CK®**

https://attack.mitre.org

**MITRE | ATLAS**

https://atlas.mitre.org/

**MITRE**

# National Cybersecurity FFRDC Research Initiative

## Goals

- Offer better understanding of the AI component

- NOT REVEAL proprietary development

- Facilitate better dialogue between vendor and potential tool adopter

- Developing a tool an organization can use

- More than cyber

**The ARCCS Framework**
**AI Relevance Competence Cost Score**

# The ARCCS Framework
# AI Relevance Competence Cost Score

*Purpose: Develop an evaluation methodology and metrics to assess the degree and effectiveness of the AI component of a commercially offered, AI-enabled product*



**Relevance**
How necessary and appropriate is the AI component?

**Competence**
How well does it do what it claims?

**Cost**
What is the cost/benefit?

MITRE

# Relevance

## Is it necessary?

## Is it central and significant?

## Is it the right tool?

- Machine learning?
- Expert system?
- The right data?

- How much functionality does AI bring?

- Sometimes, it's just gratuitous. Did we even need AI?

# Competence

## Needs Alignment



- Does it do well what you need?

## Real world demonstration



- How deal with operational issues like model drift and retraining requirements?

## Transparency



- Can we see inside the box?
- How to monitor and improve performance?

# Cost…More Than Dollars and Cents

**Tuning**
How much specialization or tuning?

**Efficiency**
Did increased accuracy slow you down?

**Security**
What vulnerabilities might be introduced?

stop sign: 99%

STOP

sports ball: 80%

STOP

*Image from: https://www.cse.gatech.edu/news/611783/erasing-stop-signs-shapeshifter-shows-self-driving-cars-can-still-be-manipulated*

# Confidence – A Modifier

- **Transparency into the model and supporting data**

- **Publications and patents**

- **White papers and publicity materials**

**MITRE**

# Strength – A Modifier Reflecting Knowns vs. Unknowns

**MITRE**

# ARCCS – What's In The Box?

**Questions**

Web browser, Spreadsheet, and questionnaire

**Scoring**

Scoring system per feature

**Inference**

Inference method on feature scores

**Guidance**

How-to guidance with guided questions and expected answers

**MITRE**

# ARCCS Browser Interface

# Coming Version

# The Promise of AI Tools



Did You Get What You Paid For?



Or Not? 😞

**MITRE**

# Take this home

- **Know that not all "AI-enabled" claims are equivalent**

- **Introduce AI-enabled tools assessments to your acquisition process**

- **Purchasers: Learn how to use and apply ARCCS**

- **Vendors: See how your product rates; use results to drive your public documentation**

*For a copy of the report/tool*
*or more information, contact: arccs@mitre.org*

*Downloads at: https://mitre.github.io/arccs/*

**MITRE**