# A New Test & Evaluation Regime for Human-AI Systems

Aditya Singh, Zoe Szajnfarber
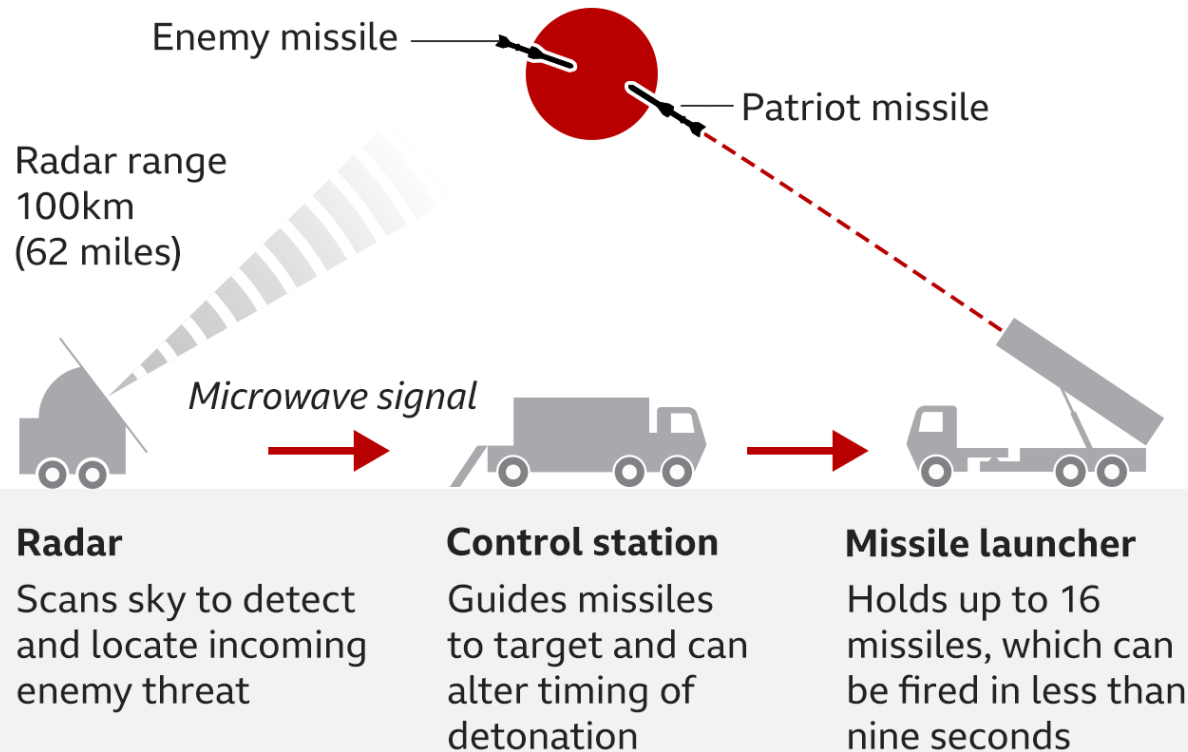
September 17, 2025

# Motivation: New Needs for T&E

**T&E has focused on the performance and reliability of the technical artifact**

**But not on how that artifact is integrated with operators, which may affect performance**
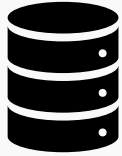


How the Patriot missile system works

Source: Raytheon Technologies



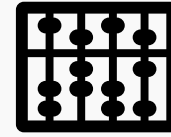Source: Netherlands Ministry of Defense

# Research Gap

Steps in AI Development & Deployment Process

**Data** → **Model Development** → **System Testing** → **Societal Impacts**

Research areas

Researchers

| | | | |
|---|---|---|---|
| ➢ Bias, poisoned data, etc. <br> ➢ Social science, CS academia | ➢ New ML Techniques <br> ➢ Frontier Labs, CS academia | ➢ Alignment, Mechanistic interpretability, etc. <br> ➢ Frontier Labs, CS academia | ➢ Trust in AI, future of work, model security <br> ➢ Humanities academia, think tanks, |

AI Integration into Human Work Systems

➢ Research often fails to consider *how* 'AI' is  integrated into workflows
➢ **How do different integrations of humans and AI change system outcomes?**

# Human-AI System Architecture is a Choice

**Human Gives Control of a Portion of Task**

Adaptative Cruise Control

**AI Takes Control of Vehicle**

Emergency Breaking Collision Avoidance

Parking Assist

**AI Assists Human**

Or Autonomous Driving Mode
(AI Drives; Human Monitors)

Lane Departure Warning

**AI Warns Human (No Action)**

Or Lane Keeping (AI Acts)

Same function can be architected in different ways

**Architecture is a decision about 1) function allocation 2) relationship b/w H&AI**

**Options are much broader than humans supervising AI or AI decision aides**

# Policy, Architecture, & Design Where is the Line?

| Policy Level | A human must have supervisory authority over any AI system's decision to use deadly force |
|---|---|

**Architecture Level**



AI Suggests a Plan of Action → Human Must Approve → AI Implements Action

Human Approver Architecture

**Design Level**



Object classified as a target
(1 min till it hits)

Fire     Do Not Fire

Bad Design



Object moving towards base at 500 mph. No friend or foe signal received. 1 min till it reaches base area.

Fire     Do Not Fire

Better Design

# Prior Work

Human Only | **Human-AI Control** | AI Only

Whose action is strictly necessary for 'the-loop' to be complete ?

## Both are Required

What is the role of the human(s)?

**Act** | **Direct AI**

Human & AI Act

AI Acts Only If Human Directs

Is the human presented one or many plans?

**Human-AI Team**

**One** | **Many**

**Human Approver** | **Human Selector**

Human-in-the-loop

## Human is Primary

Whose initiative is required for AI to act?

**AI** | **Human**

Human Acts
**AI Oversees**

Human Acts
**AI Assists**

Type of access to the control surface of system?

Type of access to the control surface of system?

**Direct** | **Indirect**

**Direct** | **Indirect**

**AI-on-the-Loop** | **AI-over-the-Loop**

**AI-along-the-Loop** | **AI-under-the-Loop**

## AI is Primary

Whose initiative is required for Human to act?

**Human** | **AI**

AI Acts
**Human Oversees**

AI Acts
**Human Assists**

Type of access to the control surface of system?

Type of access to the control surface of system?

**Direct** | **Indirect**

**Direct** | **Indirect**

**Human-on-the-Loop** | **Human-over-the-Loop**

**Human-along-the-Loop** | **Human-under-the-Loop**

When can the human(s) act?

When can the human(s) act?

Before AI Action is Final | While AI Operates

While AI Operates | After AI Acts

**Command by Veto** | **Human Supervisor**

**Executive Command** | **Human Feedback Loop**
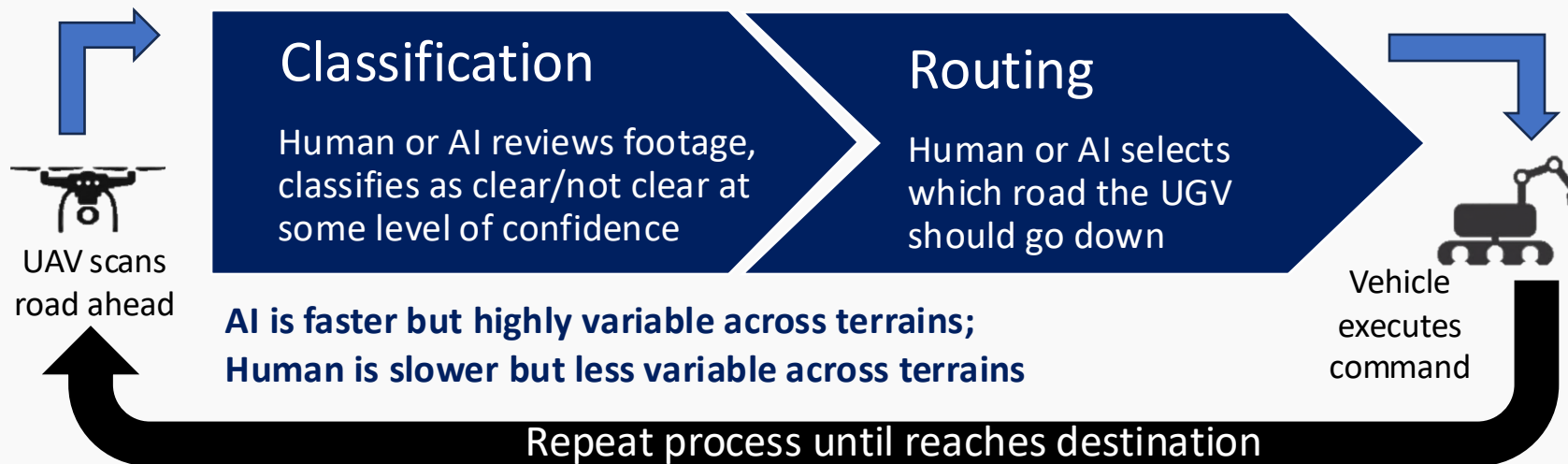
# Research Setting: Minefield Traversal

Using the framework, we modeled several architectures which determined how tasks were allocated between humans and AI and how they worked together

**Mine presence may be predicted by sending a UAV to collect data about the road.**
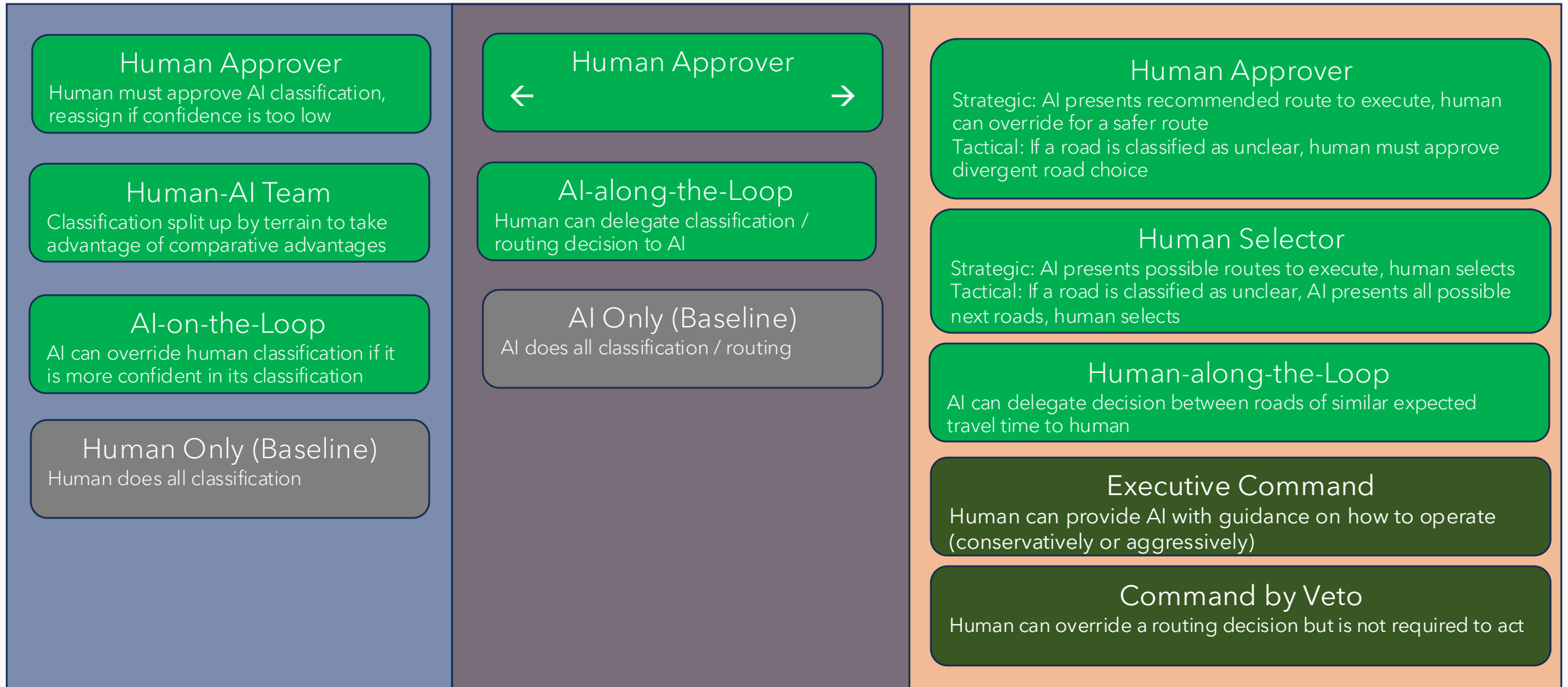
## Classification
Human or AI reviews footage, classifies as clear/not clear at some level of confidence

## Routing
Human or AI selects which road the UGV should go down

UAV scans road ahead

Vehicle executes command

**AI is faster but highly variable across terrains; Human is slower but less variable across terrains**

Repeat process until reaches destination

# Architecture Implementations in Simulation Environment

**Classification**

**Both**

**Routing**

## Classification

**Human Approver**
Human must approve AI classification, reassign if confidence is too low

**Human-AI Team**
Classification split up by terrain to take advantage of comparative advantages

**AI-on-the-Loop**
AI can override human classification if it is more confident in its classification

**Human Only (Baseline)**
Human does all classification

## Both

**Human Approver**
← →

**AI-along-the-Loop**
Human can delegate classification / routing decision to AI

**AI Only (Baseline)**
AI does all classification / routing

## Routing

**Human Approver**
Strategic: AI presents recommended route to execute, human can override for a safer route
Tactical: If a road is classified as unclear, human must approve divergent road choice

**Human Selector**
Strategic: AI presents possible routes to execute, human selects
Tactical: If a road is classified as unclear, AI presents all possible next roads, human selects

**Human-along-the-Loop**
AI can delegate decision between roads of similar expected travel time to human

**Executive Command**
Human can provide AI with guidance on how to operate (conservatively or aggressively)

**Command by Veto**
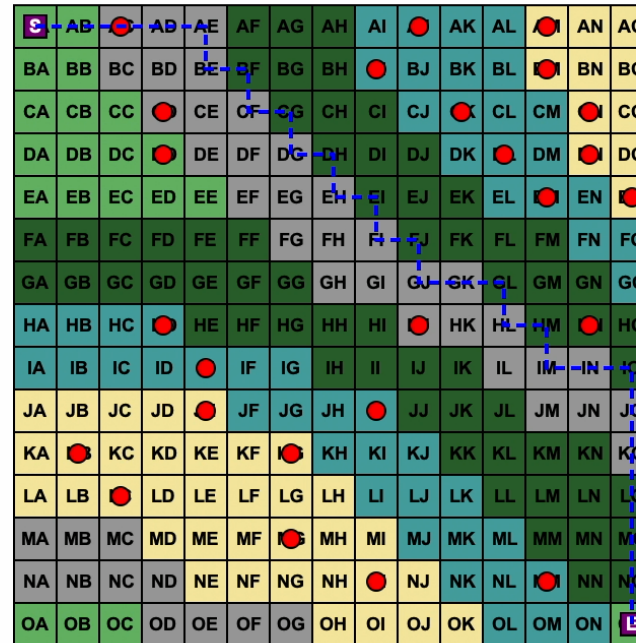Human can override a routing decision but is not required to act

# HAI Simulation Set-up



**Variables**

Treatment:
HAI Architectures

Environment:
Map size
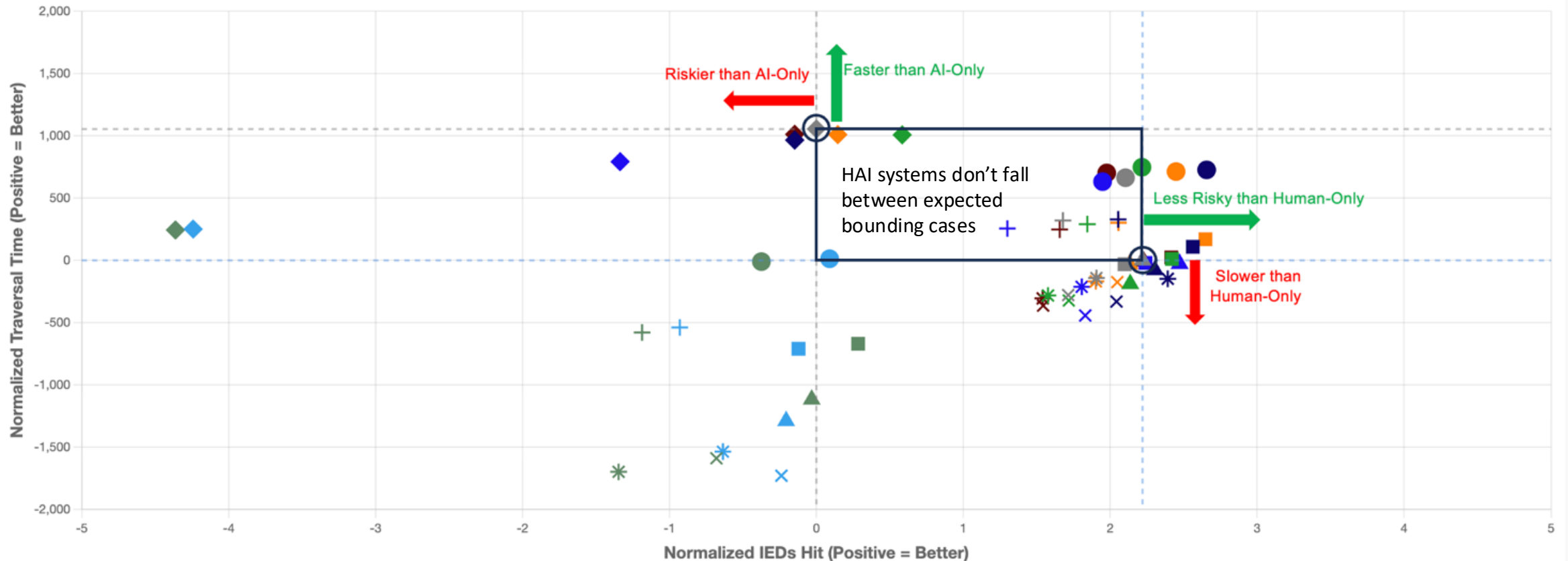Terrain
IED density
+
Human and AI
confidence

**Simulation Testbed**

**FOMs**

Performance:
Traversal time
(normalized)
Path length

Risk:
IEDs hit
(normalized)
Classification
errors
(type 1 and 2)

Legend

Swampy | Rocky | Sandy
Grassy | Wooded
Current Position | Path | IED
Start/End | Optimistic Path

Performance vs Risk (IEDs Hit): Medium Size, Heavy IED Density, Calibrated Confidence

Robustness of Results to Context

Change in Architecture Ranking Across Environmental Conditions

# Importance of Training on Interaction



Medium Size & Heavy IED Density Environment

# Findings

➤ **Architecture, Environment, and their Interaction are all significant**

  ➤ All three were statistically significant in ANOVA tests

➤ **Change in performance across environmental conditions was not uniform, consistent, or obvious**

  ➤ Seemingly innocuous changes in operating environment (increasing map size with same IED density and confidence) led to large changes in relative performance for some architectures

# Implications for Test & Evaluation

➤ Need to expand system boundary of T&E to consider human-AI architecture & interaction

　➤ Changing just *how* the human is integrated significantly changed results while holding the technical performance constant

➤ Human-AI systems testbeds can:

　➤ Reveal non-obvious tradeoffs and interactions

　➤ Understand how changing variables affect system outcomes

　➤ Identify which architectures that are robust / sensitive to expected operating environment

# Thank You

asingh25@gwu.edu

# Classification Architectures

## Human-AI Team

Leverages complementary strengths based on historical performance

UAV Scans Road → Check Terrain →

- Rocky/Grassy Terrain → Human Classifies
- Swampy/Sandy/Wooded → AI Classifies

→ Final Classification

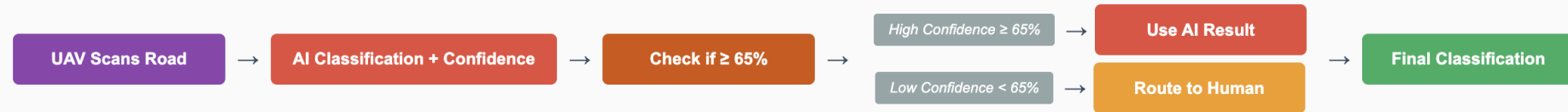## Human Approver

Threshold adjusts dynamically: decreases per correct AI decision (building trust), increases 20% after IED encounter (betrayal).

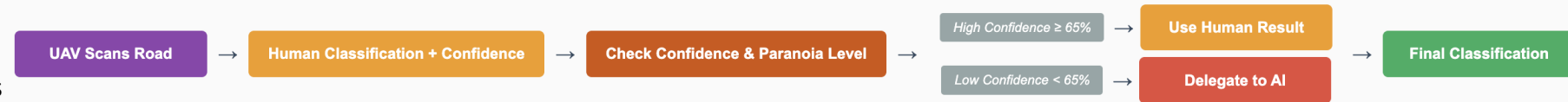UAV Scans Road → AI Classification + Confidence → Compare to Human Threshold →

- Confidence ≥ Threshold → Accept AI Result
- Confidence < Threshold → Human Reclassifies

→ Final Classification

## Human-along-the-Loop

Simple rules-based system with fixed threshold, reassigns low-confidence cases to human expert.

UAV Scans Road → AI Classification + Confidence → Check if ≥ 65% →

- High Confidence ≥ 65% → Use AI Result
- Low Confidence < 65% → Route to Human

→ Final Classification

## AI-along-the-Loop

Models paranoia that increases with consecutive "clear" classifications. Paranoia resets when mines are found or roads marked unclear.

UAV Scans Road → Human Classification + Confidence → Check Confidence & Paranoia Level →

- High Confidence ≥ 65% → Use Human Result
- Low Confidence < 65% → Delegate to AI

→ Final Classification

## AI-on-the-Loop

AI monitors and only overrides when it disagrees AND has significantly higher confidence (+10% margin).

UAV Scans Road → Human Classification / AI Classification (Parallel) → Compare Results & Confidence →

- Agreement OR AI not significantly better → Use Human Result
- Disagreement + AI confidence >Human +10% → AI Override
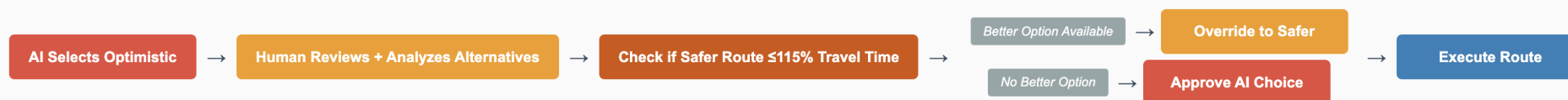
→ Final Classification

# Routing Architectures (Strategic vs Tactical)

**Human makes decisions about what type of route should be executed**

## Human Approver (Strategic)

Human can substitute a safer path (more favorable terrain) if expected travel time is within 15% of shortest path.

| AI Selects Optimistic | → | Human Reviews + Analyzes Alternatives | → | Check if Safer Route ≤115% Travel Time | → | *Better Option Available* → Override to Safer | → | Execute Route |
| | | | | | | *No Better Option* → Approve AI Choice | | |

## Human Selector (Strategic)

**Route Options:** 1) Shortest path, 2) AI-favorable, 3) Human-favorable.
**Selection Logic:** Human-AI Team → shortest path. AI-dominant systems → AI-favorable. Human-dominant systems → Human-favorable.

| Generate 3 Routes | → | Human Analyzes Options | → | Select Best Route | → | Execute Route |

**Human makes decisions when an issue occurs**

## Human Approver (Tactical)

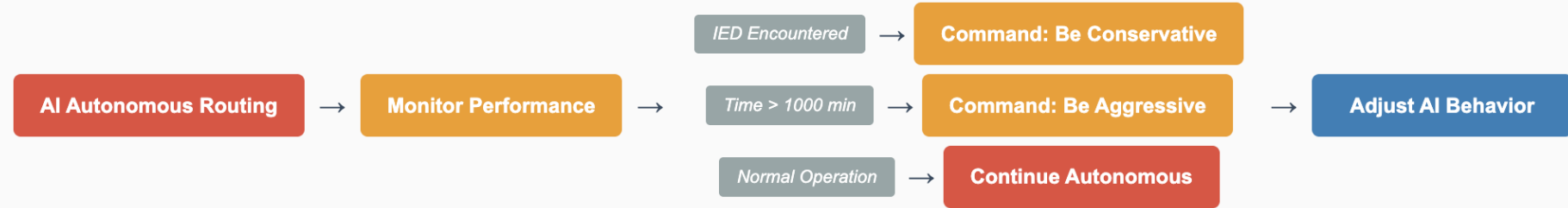If confidence is low, human selects next appropriate road classified as clear with highest confidence.

| *If Blocked* | → | Scan All Roads | → | AI Suggests Lowest Expected Time | → | Check AI Confidence vs Threshold | → | *High Confidence* → Accept Suggestion | → | Select Road |
| | | | | | | | | *Low Confidence* → Override Choice | | |

## Human Selector (Tactical)

Human weighs expected travel time, progress toward goal, and AI classification performance in different terrains.

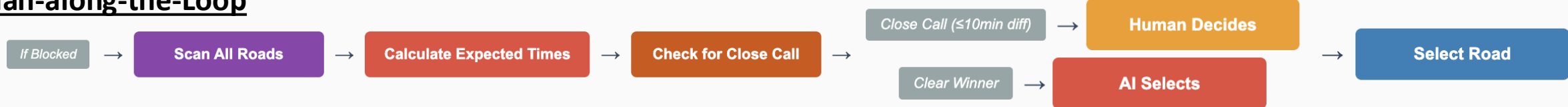| *If Blocked* | → | Scan All Node Roads | → | Multi-Criteria Analysis | → | Select Road |

# Routing Architectures (Only One Implementation)

## Executive Command

Human adjusts AI behavior. Conservative mode prioritizes roads where AI has high classification confidence over pure expected value. Aggressive mode emphasizes progress toward end node over safety margins
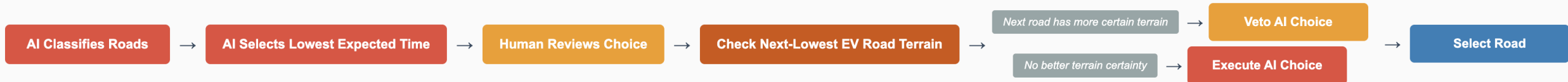
AI Autonomous Routing → Monitor Performance →

IED Encountered → Command: Be Conservative

Time > 1000 min → Command: Be Aggressive → Adjust AI Behavior

Normal Operation → Continue Autonomous

## Human-along-the-Loop

If Blocked → Scan All Roads → Calculate Expected Times → Check for Close Call →

Close Call (≤10min diff) → Human Decides

Clear Winner → AI Selects

→ Select Road

When two or more road options have expected travel times within 10 minutes, AI delegates to human expertise. Human Selects the road classified as clear with highest confidence level that makes progress toward destination.

## Command by Veto

AI Classifies Roads → AI Selects Lowest Expected Time → Human Reviews Choice → Check Next-Lowest EV Road Terrain →

Next road has more certain terrain → Veto AI Choice

No better terrain certainty → Execute AI Choice

→ Select Road

Human can reject AI's lowest expected travel time choice if the next-lowest EV road has more certain terrain type for AI classification but is not required to act for AI to operate