

AI Support Systems in the Military

Division of Science and Mathematics
University of Minnesota, Morris

John Gulon



April 9, 2026



Road Map

- Why AI Decision Support Matters
- What is ReLeGSim?
- How the Model Architecture works
- How Reward & Training shape Behavior
- Experimental Results

Why is AI in the Military?

- Data rich environment
- Human limits (OODA: Observe, Orient, Decide, Act)
- Faster, Adaptive Decision Support
- Decision Support, **not** Replacement

What is ReLeGSim?

- Reinforcement Learning-based Ground Simulation
- Battalion-level Simulation (BLUE vs RED)
- Grid-based Battlefield
- AI Generates Tactical Commands
- Reinforcement Learning
- Designed by Möbius et al.

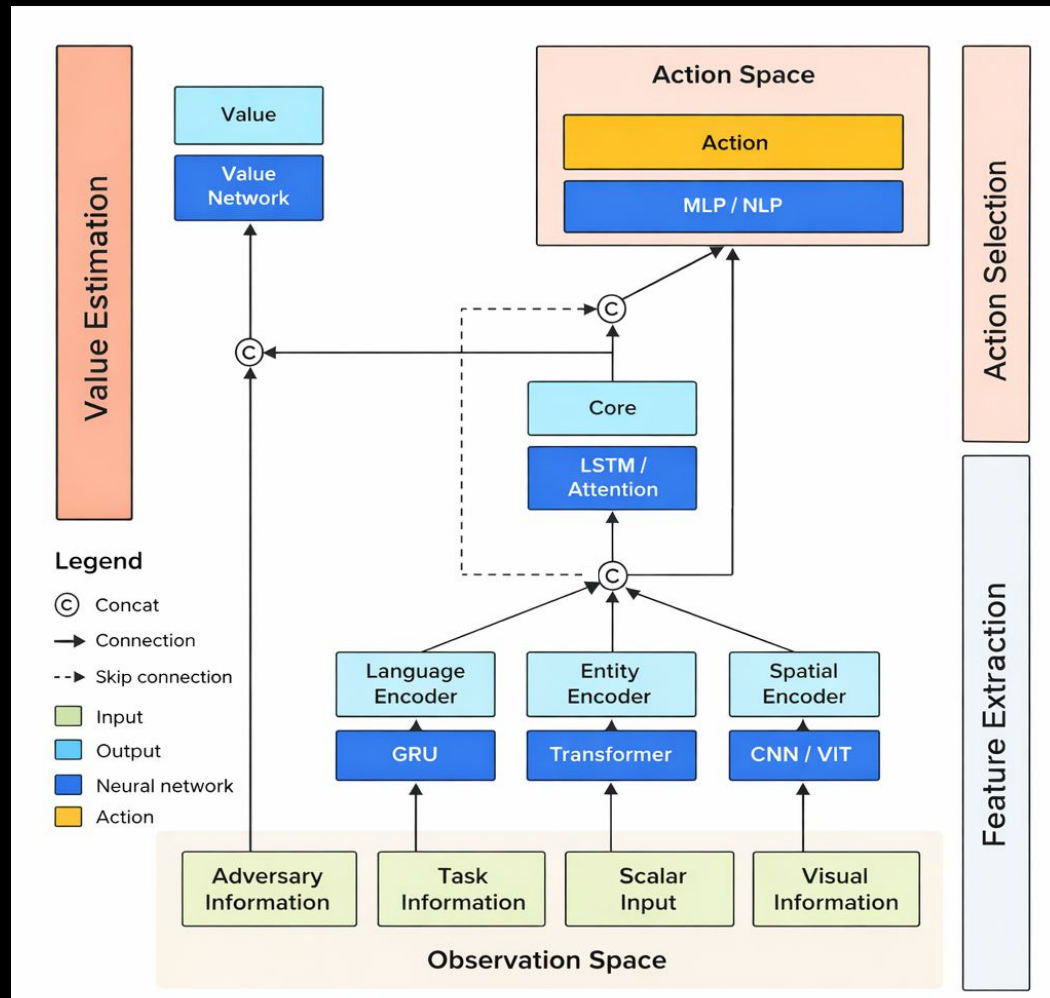
Map adapted from Möbius et al. 2024



Model Architecture

- Multi-modal input
- Feature Extraction
- Decision Network
- Action & Value Outputs

Disclaimer: Equations used are theoretical, they may or may not be used by sources.



Reward Function

Measuring Long-term Success

$$R(s_t, a_t) = w_1 r_{\text{mission}} + w_2 r_{\text{enemy}} + w_3 r_{\text{friendly}}$$

Define what the AI considers “good behavior” (Policy)

s_t = state, time step.

a_t = action, time step.

Why Reward Design Matters

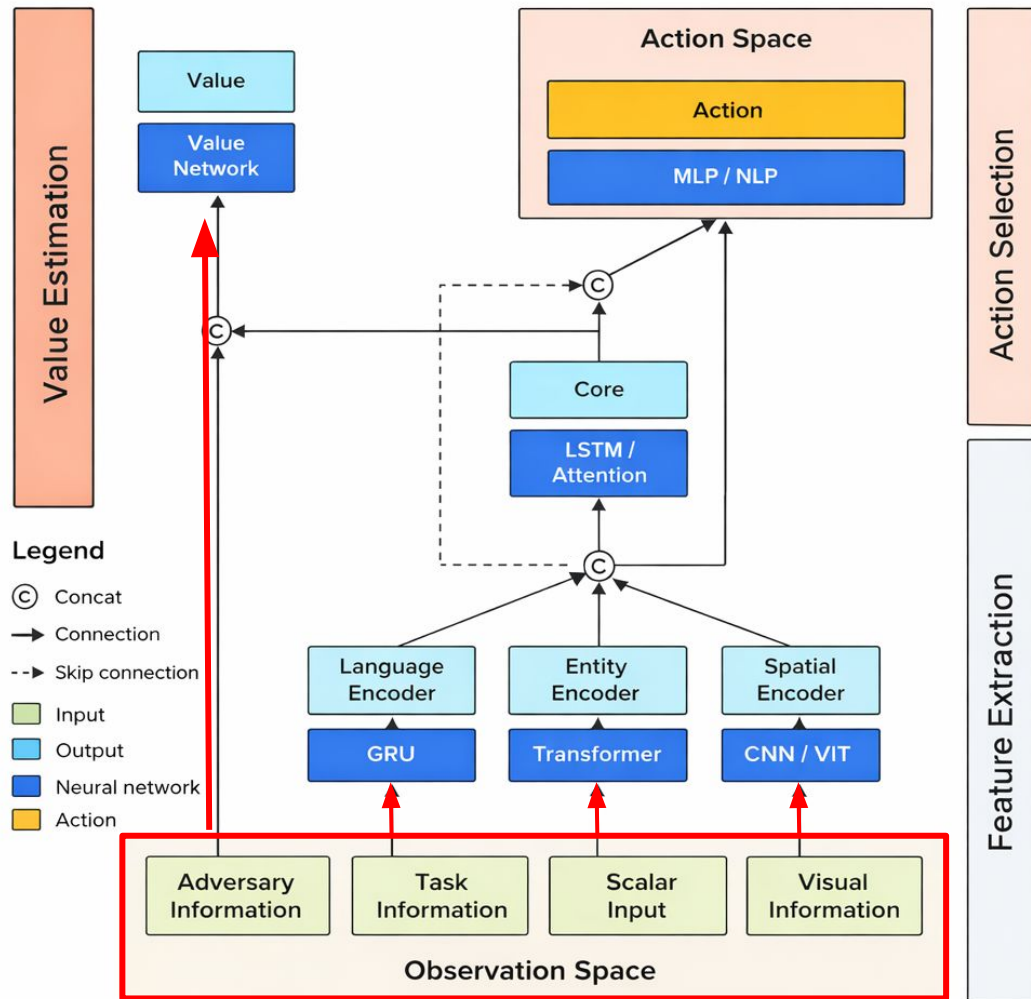
- Reward defines Behavior

Simple Rewards = Narrow Strategies

Multi-objective Rewards = Realistic Tactics

- Doctrine Alignment

Observation Space



Observation Space

- Adversary Information
- Task Information
- Scalar Input
- Visual Information

$$S_t = f(\text{enemy}, \text{task}, \text{friendly}, \text{terrain})$$



Equations Road Map

1. Input State

- Enemy armored units near B4
- Task: Capture Objective, Minimize Losses
- Friendly companies w/Artillery Support
- Forest, road, open terrain (Map)

$$S_t = f(\text{enemy, task, friendly, terrain})$$

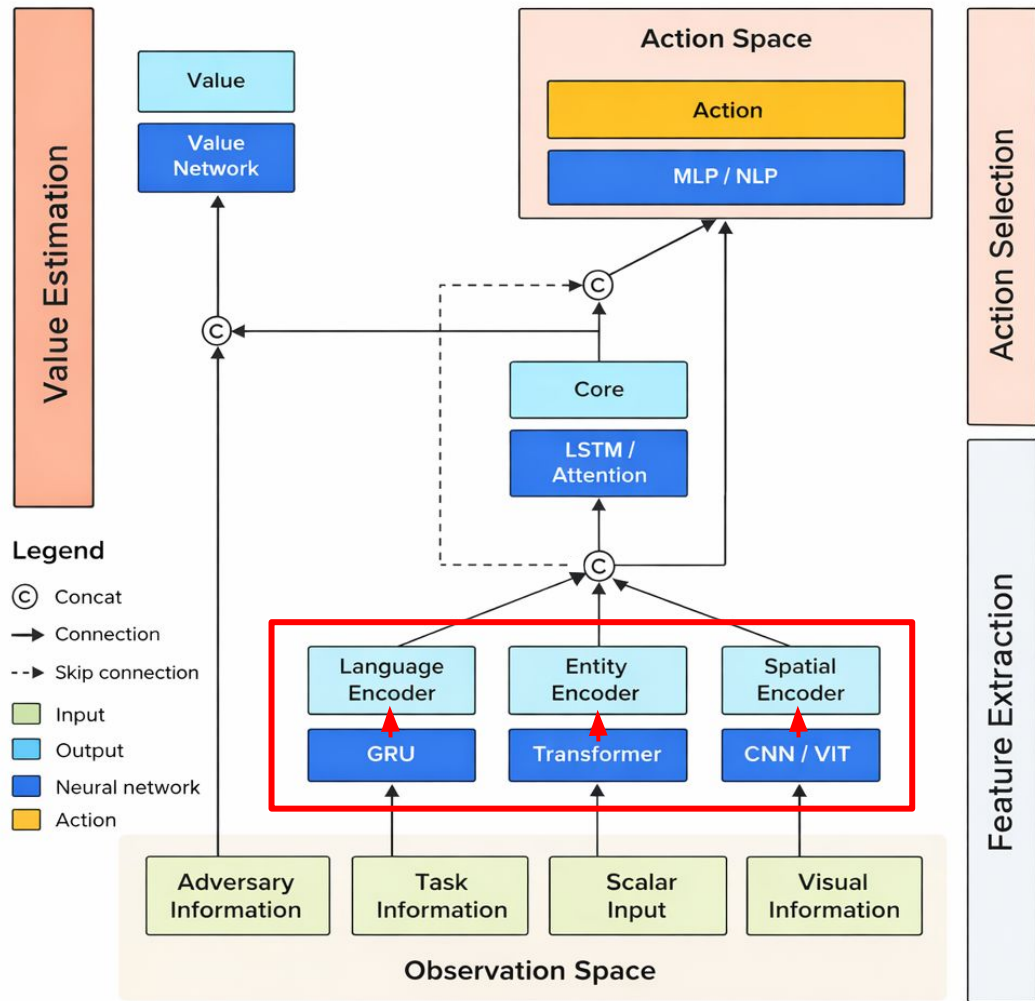
Observation (Model View)

- enemy_forces_table
- own_forces_table
- combat_support_table
- objective
- image
- prev_action
- action_mask (valid actions)

Feature Extraction

- GRU -> Command Intent
- Transformer -> Unit
- Relationship
- CNN/ViT -> Spatial Features

$$h_t = [h^{\text{lang}}, h^{\text{entity}}, h^{\text{visual}}, h^{\text{scalar}}]$$



Core Network

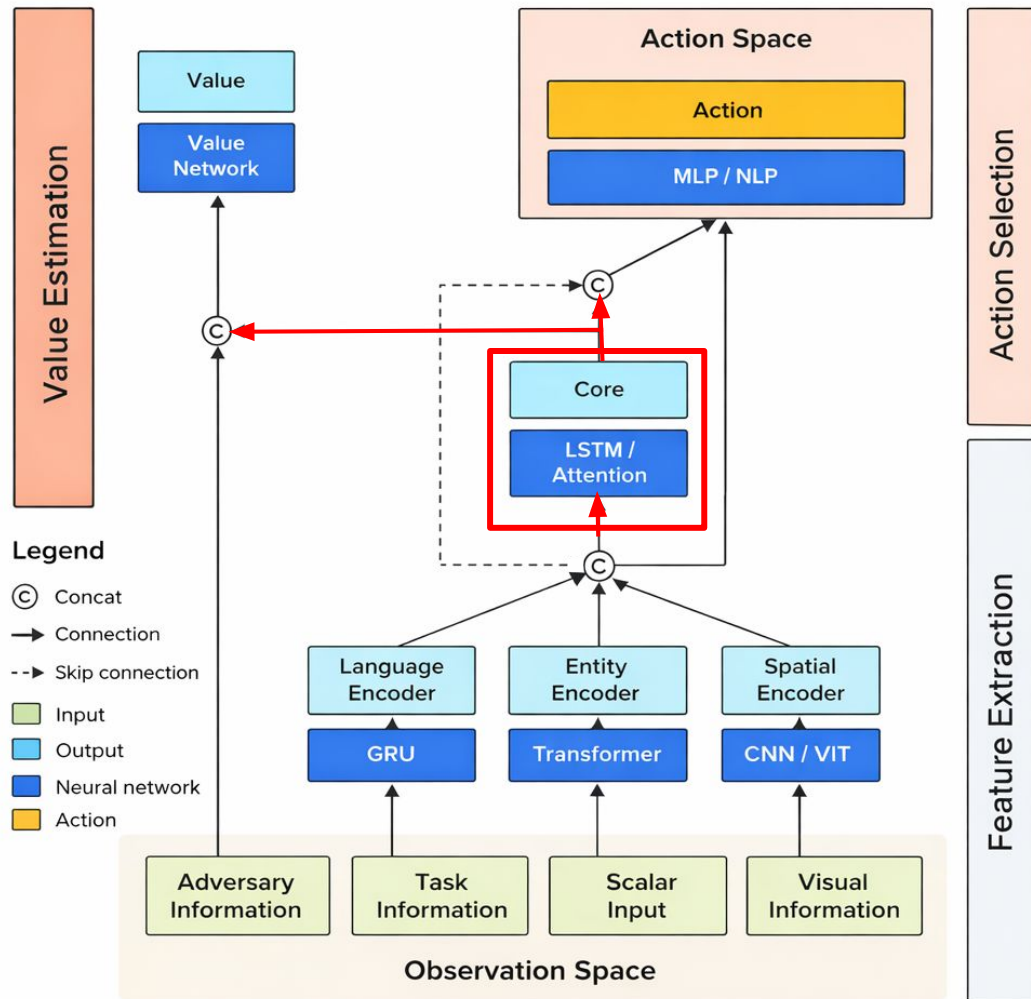
- LSTM -> Memory Over Time

$$h_t = \text{LSTM}(h_t, h_{t-1})$$

- Attention -> Focus on

Important Info

$$c_t = \sum \alpha_i x_i$$



Equations Road Map

2. Core

$$h_t = [h^{\text{lang}}, h^{\text{entity}}, h^{\text{visual}}, h^{\text{scalar}}]$$

x_i = individual elements/features inside h_t
($h_t = [x_1, x_2, x_3, \dots]$)

Core (Recurrent + Attention)

$$h_t = \text{LSTM}(h_t, h_{t-1})$$

$$c_t = \sum \alpha_i x_i$$

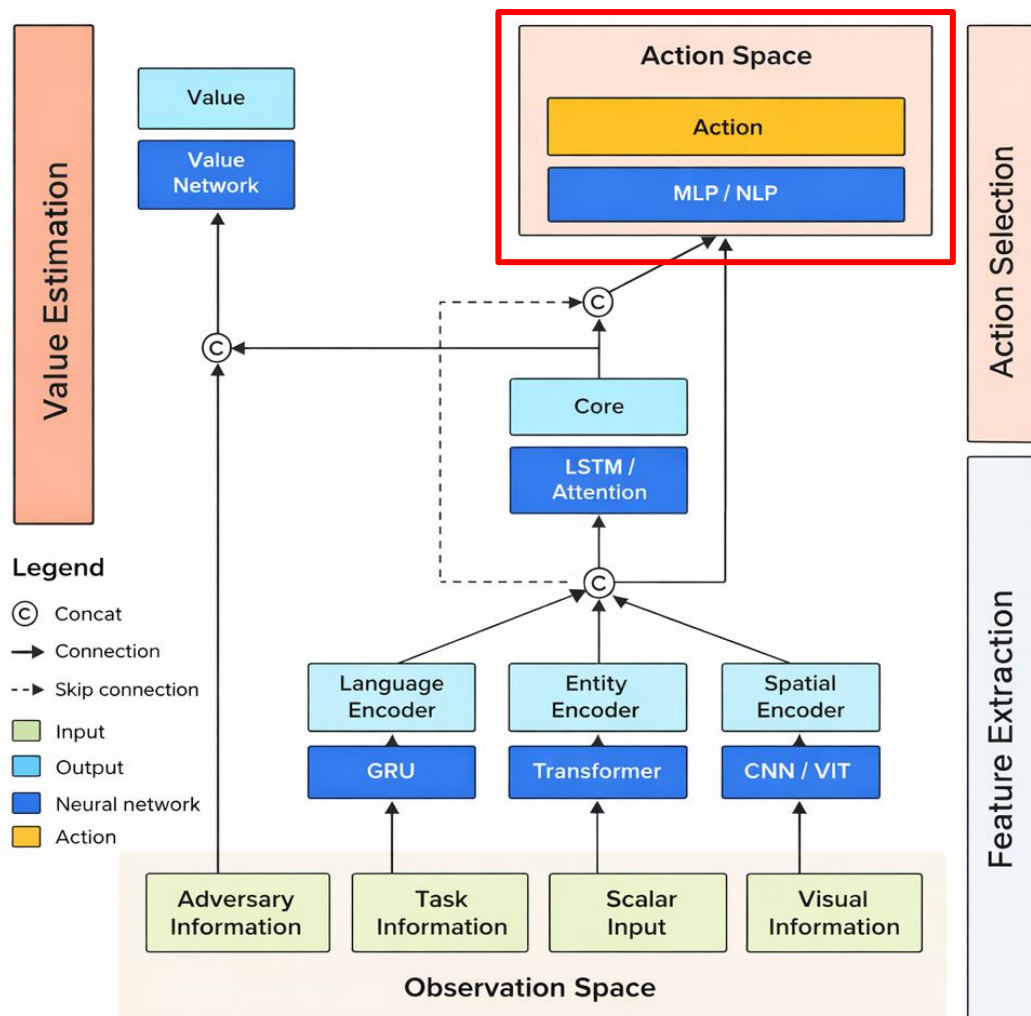
α_i = learned importance weights

Action Selection

Outputting Commands

- Move
- Attack
- Observe

$$\pi(a \mid s) = P(a_t \mid s_t)$$
$$a_t = \operatorname{argmax} \pi(a \mid s)$$



Equations Road Map

3. Action Selection

h_t^{core}

Policy (Probability Distribution)

$$\pi(\mathbf{a} \mid \mathbf{s}) = \mathbf{P}(\mathbf{a}_t \mid \mathbf{s}_t)$$

- Maps State \rightarrow Probabilities over actions
- Constrained by:
 - Action mask
 - Available resources

Action Choice

$$\mathbf{a}_t = \mathbf{argmax} \pi(\mathbf{a} \mid \mathbf{s})$$

- Selects highest probability action
- Outputs command
 - Move
 - Attack
 - Observe

Equations Road Map

4. Value Estimation

h_t^{core}

Value Function

$$V(s_t) = E [\sum \gamma^k R]$$

- Estimates total future reward
- Evaluates current state quality

Output: $V(s_t) \in \mathbb{R}$, a single scalar value for state quality.

Full Learning Objective

Maximizing **Long-term** reward while following mission **priorities**.

$$J(\pi) = E [\sum \gamma^t R(s_t, a_t)]$$

Combining action, value, reward into a single objective.

Training Loop

- observe state s_t
- compute features and core representation
- output action and value
- execute action in simulator
- receive reward and next state
- update policy/value network
- repeat over many time steps and episodes

Reward Function Reminder

Measuring Long-term Success

$$R(s_t, a_t) = w_1 r_{\text{mission}} + w_2 r_{\text{enemy}} + w_3 r_{\text{friendly}}$$

Define what the AI considers "good behavior" (Policy)

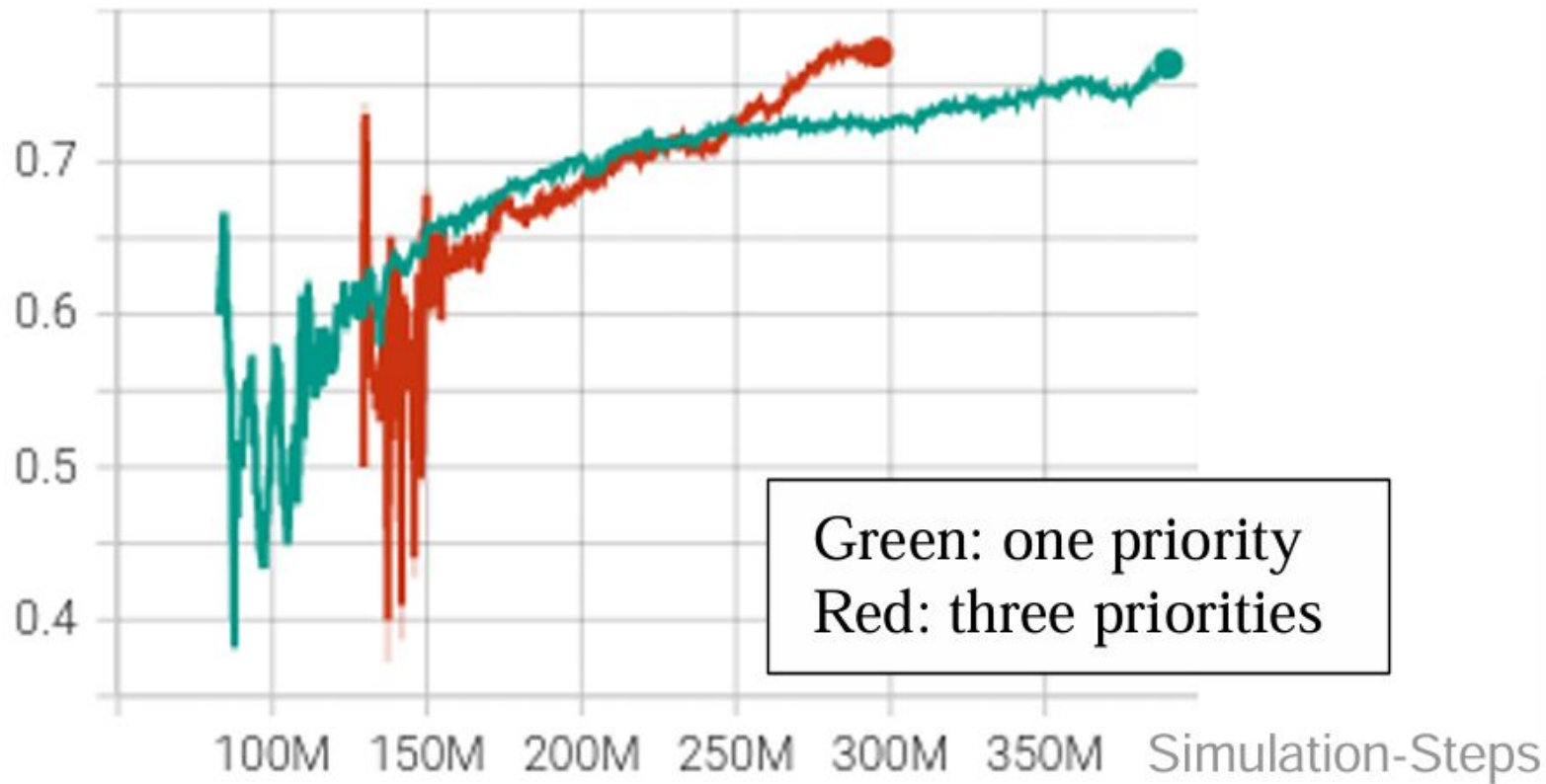
s_t = state, time step.

a_t = action, time step.

Experimental Results

Training Results from Möbius et al. 2024

Win rate in %



Experimental Results

Early Training

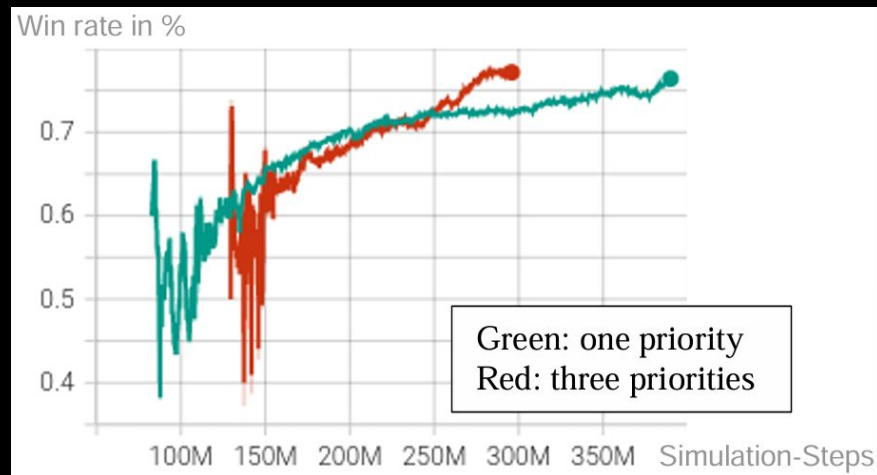
- Multi-objective = slower

Later Training

- Multi-objective = better performance

Reason


- More exploration
- More realistic strategies



Connection to Doctrinal Alignment

Results Support the Doctrine

- Real Military Decisions = Multi-Objective
- Reward = Doctrine Encoding
- Better Performance = Better Alignment



Human-in-the-Loop

- AI Suggests
- Human Decides
- Prevents Misalignment

Key Takeaways

- Multi-modal AI = Battlefield Awareness
- Reward = Doctrine
- Architecture = Decision Pipeline
- Risk = Misalignment
- Human Control = Priority



Questions?

Citations

Main Source

Michael Möbius, Daniel Kallfass, Matthias Flock, Thomas Doll, and Dietmar Kunde. 2024. Incorporation of Military Doctrines and Objectives into an AI Agent via Natural Language and Reward in Reinforcement Learning. In Proceedings of the Winter Simulation Conference (WSC '23). IEEE Press, 2357-2378.

Supporting Source

Michael Möbius, Daniel Kallfass, Thomas Doll, and Dietmar Kunde. 2023. AI-Based Military Decision Support Using Natural Language. In Proceedings of the Winter Simulation Conference (WSC '22). IEEE Press, 2082-2093.

Supporting Source