

Australian Government Department of Defence

# Machine-Learning & Recommender Systems for C2 of Autonomous Vehicles

### Glennn Moy

on behalf of

Don Gossink, Glennn Moy, Darren Williams, Kate Noack Josh Broadway, Jan Richter, Steve Wark

Planning and Logistics, Decision Sciences, DST Group, Australia

Science and Technology for Safeguarding Australia

## **Human-Autonomy Teaming**

Overall project – Balances two components:

### Autonomy Lead: Don Gossink



HAT - Autonomy Team:

- Don Gossink
- Glennn Moy
- Darren Williams
- Katherine Noack
- Josh Broadway
- Jan Richter
- Steve Wark

....

:::

....

### **HAT Context**

**Recommendations for Command & Control of** Multiple UxVs



....

.

÷

÷-... **...** 

.

÷

**.**... ....

....

•

. ....

### IMPACT

- U.S. Project aimed at:
  - Developing Intelligent Multi-UxV Planner with Adaptive Collaborative/Control Technologies
  - Key Concept: High-level, goal-oriented plays



Science and Technology for Safeguarding Australia

### **Recommender System & Play-Monitor:**

- Goal:
  - Develop advanced Recommender System(s) to reduce the cognitive burden on operators through:
    - Recommendations, alerts and constraints. over the top of

"Human on the loop"

Lower-Level Autonomy



÷

## **Top-Level Architecture:**



### **HAT Challenge**

- Limited access to IMPACT or other multi-UxV control system
- Need to integrate recommendations into various (unknown) system components
  - Solution:
    - Recommender *loosely coupled* to the underlying system/simulation.

- Recommender that can *learn* recommendations at a range of C2 levels.
- Recommender techniques that work:
  - When heuristics are not known

•

In new contexts (not previously prepared)

## **Inspiration & Approach**

.....

•

Build something that can learn... (like we learn???)





(NB: about gathering and structuring data to learn...not how we learn)

### **Techniques & Requirements:**

ŀ

....

**...** 

.....

•

|                                 | Self Learning  | Learning from Others   |
|---------------------------------|--|--|
| Experience                      | <ul> <li>Trying:</li> <li>Techniques: Reinforcement Learning</li> <li>Requirements: A simple model that machine can control that it can undertake re-enforcement learning on.</li> </ul>                                       | <ul> <li>Watching:</li> <li>Techniques: Supervised Deep Learning</li> <li>Requirements: Data from expert playing a simulation that it can watch.</li> <li>Human Expert or Heuristic player that mimics an expert</li> </ul>  |
| Logic,<br>Rules &<br>Heuristics | <ul> <li>Reasoning:</li> <li>Techniques: Monte-Carlo State<br/>Exploration/Search, Logic, Planning.</li> <li>Requirements: Implement efficient<br/>heuristic/search algorithms for<br/>exploring large state-space.</li> </ul> | <ul> <li>Explaining:</li> <li>Techniques: Heuristics, Agent (BDI),<br/>Algorithms, Math/Planning, Abductive<br/>logic</li> <li>Requirements: <ul> <li>A language for expressing heuristic<br/>rules.</li> <li>Logic for constraining solutions.</li> </ul> </li> </ul> |

### **Architecture:**



## **Recommendation Hierarchies**

- Recommender Agents:
  - Implemented Recursively:

÷

- Hierarchy of recommendation agents.
- Key Concepts:
  - Decomposing / Triggered Recommender Agents Elastic autonomy
  - Executable vs Non-Executable Recommendations
- Simulation:

11

Publishes its own capabilities for accepting/executing recommendations.



#### Science and Technology for Safeguarding Australia

## **1. Simulation**



- Initial UAV Control Simulations
  - Modular/Plug-and-play
    - Ultimately to be replaced by IMPACT
  - Fast

13

In a training mode for rapid learning

....

÷

....

•

÷

.....

Machine or Human controllable





| QUIT Start Process Stop Process Switch Directions! Unknown Action | Í | <b>%</b> tk |               | - • •        |                    |                |
|---|---|-------------|---------------|--------------|--------------------|----------------|
|   |   | QUIT        | Start Process | Stop Process | Switch Directions! | Unknown Action |

### **2. First Heuristic Recommenders:**

- Initial low-level *Executable* Recommenders
- Heuristics & Play Implementations
  - (a) Search for detection
    - Air-Expanding Square at/on a Point
    - Air Sector Search
    - Air Inspect Point
  - (b) Track Detections







### **3. First Reinforcement Recommender**

- Initial Reinforcement Recommender
  - Deep Q Reinforcement Learning :
    - Combines reinforcement learning with deep neural network.
    - Originally used to play Atari Games DeepMind/Google

Reinforcement Learning Challenges:

....

- Credit-assignment problem
- Explore-exploit dilemma





### **Deep Q Learning**

#### **Credit-Assignment Problem:**

• Q(s, a) represents the maximum discounted future reward when we perform action a in state s (and get to state s' with reward r.)

- 
$$Q(s, a) = R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q(s', a')$$
 - Bellman Eqn

- 
$$Q(s, a) = r + \gamma \max_{a'} Q(s', a')$$

- Iterative Updates (Training):
  - **Prediction**: Q(s, a)
  - **Target**:  $r + \gamma \max_{a'} Q(s', a')$

$$- L = \frac{1}{2} \left[ r + \gamma \max_{a'} Q(s', a') - Q(s, a) \right]^2$$

• Experience Replay:

**...** 

16

Store all experiences <s, a, r, s'> in memory.

- Train on random mini-batches
- Propagates rewards back in time.

### **Deep Q Learning**

### **Explore/Exploit Dilemma**

- $\varepsilon$  greedy exploration
  - With probability  $\varepsilon$  choose a random action, otherwise go with highest Q value action.
  - $\varepsilon$  starts at 1.0 (always random) and slowly decreases (eg to 0.1 mostly policy-based choices)



– Initial Task:



18

**.**....

## Example 1:

Actions: Up, Down, Left, Right (Constant Velocity, No Pause) Reward: -(Distance from Goal)



....

÷

÷

÷

## Example 2:

Actions: Rotate Left, Right, Speed up, Slow Down, Do Nothing Reward: -(Distance from Goal)

### **Deep Q Reinforcement Learning** (UAV learning to follow a boat)

### **Available Control Actions:**

- Rotate Left
- Rotate Right
- Speed up
- Slow Down

**∷**•

Maintain Constant Velocity

#### **Reward:**

20

Negative-distance from Goal (Boat)

## **Current/Ongoing R&D:**

General:

21

- Implement Red "avoid detection" behaviour & blue agent vs red agent play
- Learn/iterate to higher-level strategies / C2.

### Deep Q R&D:

- Alternative Neural Architectures:
  - Impact on learning rate
  - Scalability with action-complexity
- Multi-Agent Learning:
  - Single control agent vs multiple learning agents vs hybrid
- Human-Guided Learning
  - Learning from human interactions as well as self-generated.

#### Other Recommender Agents:

- Agent Hierarchies
  - More complex (hierarchical) reward functions and interactions between recommender agents
- Bayesian Inference for threat heat-maps

### Questions

