

# ACAS Xu Policies via Deep Reinforcement Learning

KYLE JULIAN

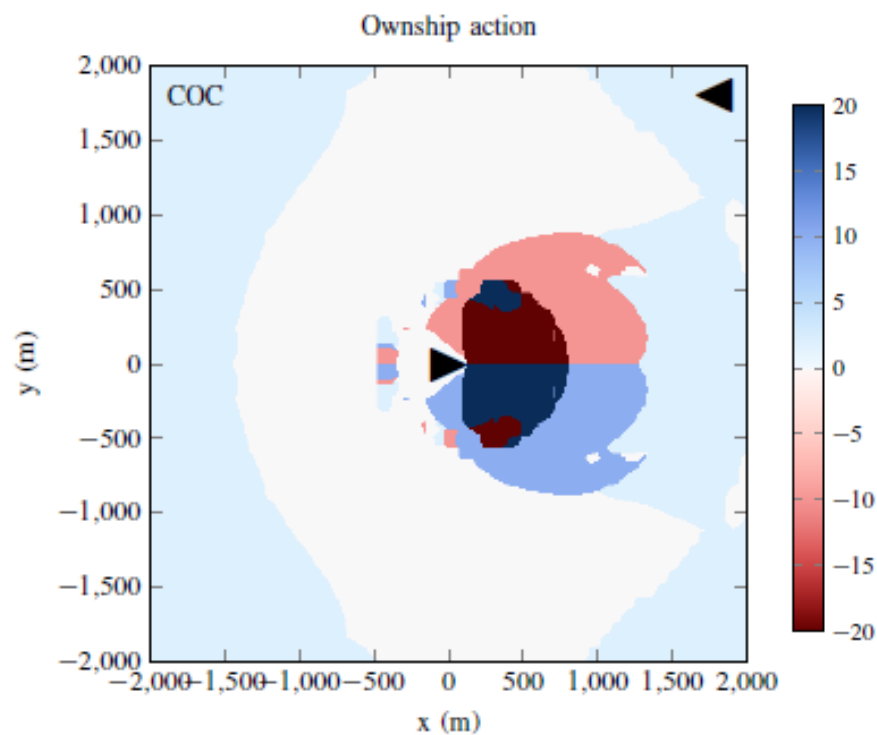
March 9, 2015

# Overview

- Prior Work
- Problem Task and Goals
- Approach
- Results
- Next Steps

## Prior Work

- Airborne Collision Avoidance System X for UAS (ACAS Xu)
- Formulated as POMDP
- Pairwise encounters
- Solved with QMDP



## Project Motivation

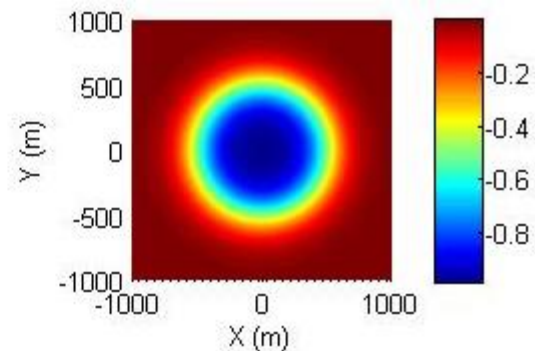
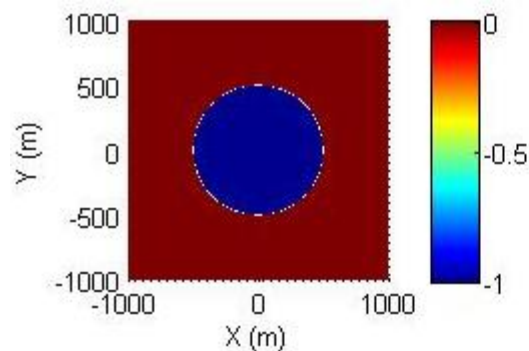
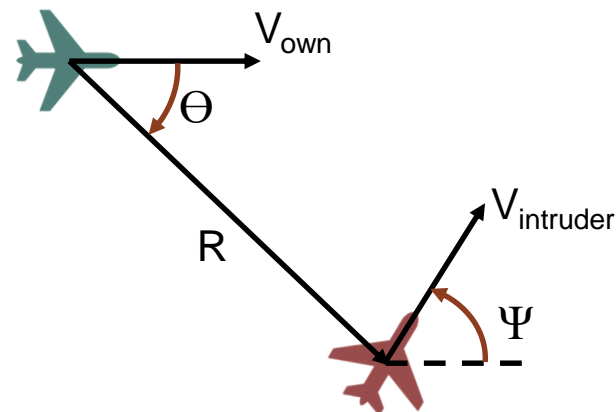
- Q table has 14 million state-action pairs
  - Future tables have up to 600 million Q values
  - Requires large space in memory
- Large policies take a full day to compute
  - Too long to study effect of reward structure on policies

## Project Goals

- Solve problem through deep reinforcement algorithm
- Compress representation of Q values
  - Estimated through network parameters
- Compute network quickly

# Problem Formulation

- State:
  - $R: [0, 3000\text{m}]$
  - $\Theta, \Psi: [-\pi, \pi]$
  - $V_{\text{own}}, V_{\text{intruder}}: [10 \frac{m}{s}, 20 \frac{m}{s}]$
- Actions (bank angles  $\phi$ ):
  - $\{-20^\circ, -10^\circ, 0^\circ, 10^\circ, 20^\circ, \text{COC}\}$
- Rewards, four possible penalties:
  - $-\frac{1}{1.0 + e^{-(r_{\text{min}} - r) * C}}$
  - $-0.0002 * \phi^2$
  - $-0.03$  if action  $\neq \text{COC}$



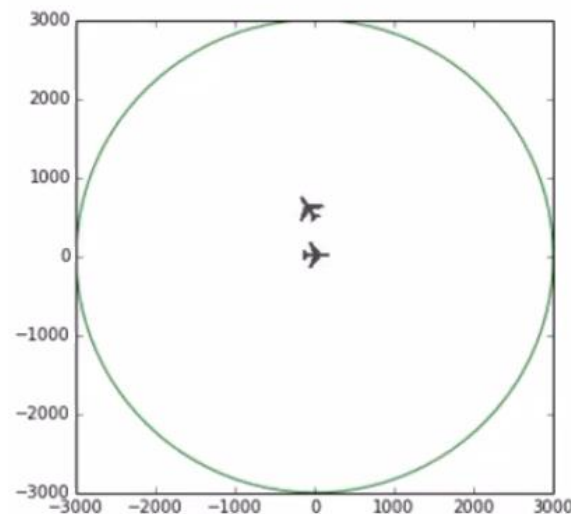
# Dynamics

- Bank angle and speeds are noisy measurements
- True values are drawn from normal distributions:

- $$\phi = \begin{cases} N(action, \sigma = 4^\circ), & action \neq COC \\ N(0^\circ, \sigma = 10^\circ), & action = COC \end{cases}$$

- $$V = N\left(V_{meas}, \sigma = 2 \frac{m}{s}\right)$$

- $$x_{i+1} = r_i \cos(\theta) + (v_{intruder} \cos(\psi) - v_{own}) * dt$$
- $$y_{i+1} = r_i \sin(\theta) + v_{intruder} \sin(\psi) * dt$$
- $$\Delta\psi = -g * \frac{\tan(\phi)}{v_{own}} * dt$$
- $$r_{i+1} = \|[x_{i+1}, y_{i+1}]\|_2$$
- $$\theta_{i+1} = atan2(y_{i+1}, x_{i+1}) + \Delta\psi$$
- $$\psi_{i+1} = \psi_i + \Delta\psi$$



Time: 0.2  
Speed Own: 20.00  
Speed Intruder: 13.78  
Intruder Heading: 123.4  
Range: 608  
Bank Commanded: 1  
True Bank: -2.78

## Techniques Applied to DQNs

- Stochastic transitions
  - Sampled  $K$  rewards and next states
  - Used average rewards and  $Q$  values
- Smoother Rewards
- Stored multiple  $(S, A, R, S')$  tuples between network updates
  - Larger batch sizes



# Pseudocode

---

**Algorithm 1** ACAS Xu DQN Algorithm

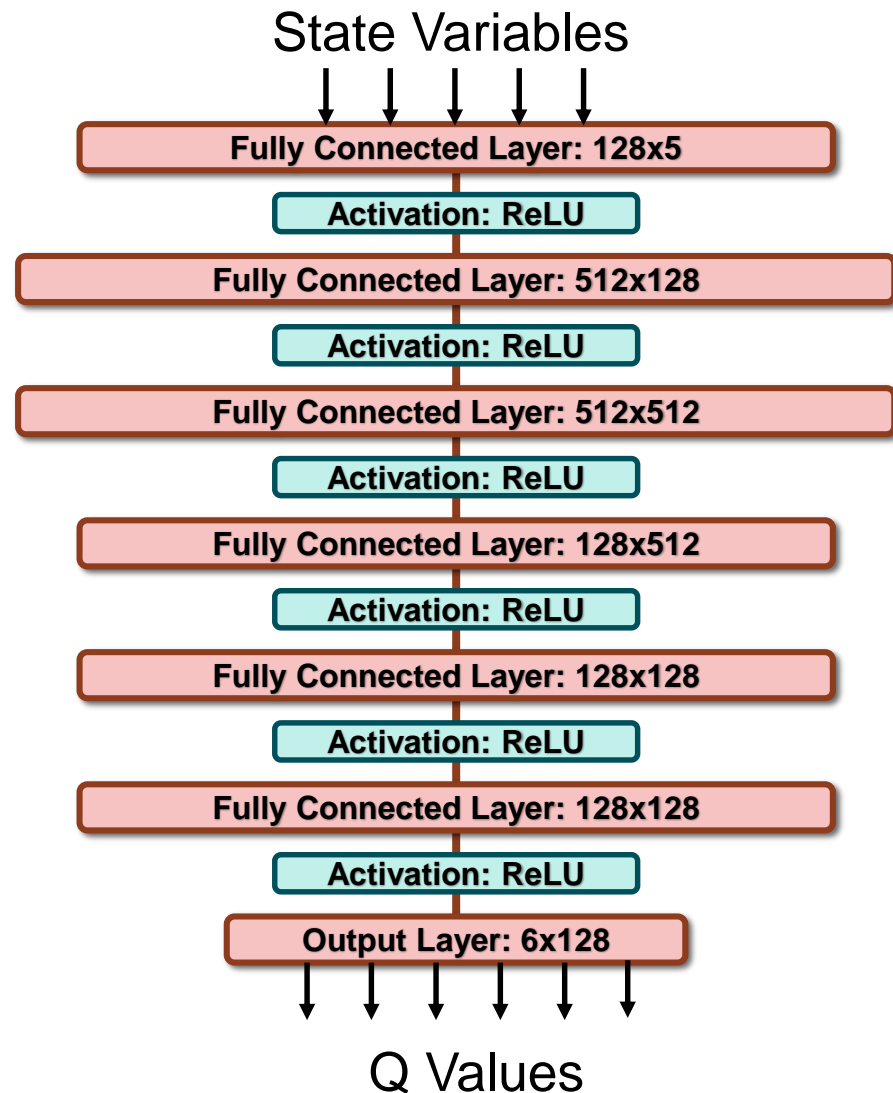
---

```
1: initialize DQN
2: initialize replayMem
3: initialize stateGen
4: repeat
5:   for  $i \leftarrow 1, K$  do
6:      $action \leftarrow DQN.getAction(state)$ 
7:      $nextStates, rewards \leftarrow stateGen.nextState(state, action)$ 
8:      $replayMem.store(state, action, rewards, nextStates)$ 
9:      $state \leftarrow nextStates[0]$ 
10:    if  $state.Range > 3000$  then
11:       $state \leftarrow stateGen.randomState()$ 
12:    if  $repMem.canTrain()$  then
13:       $batch \leftarrow repMem.sampleBatch()$ 
14:       $DQN.train(batch)$ 
15: until converged
```

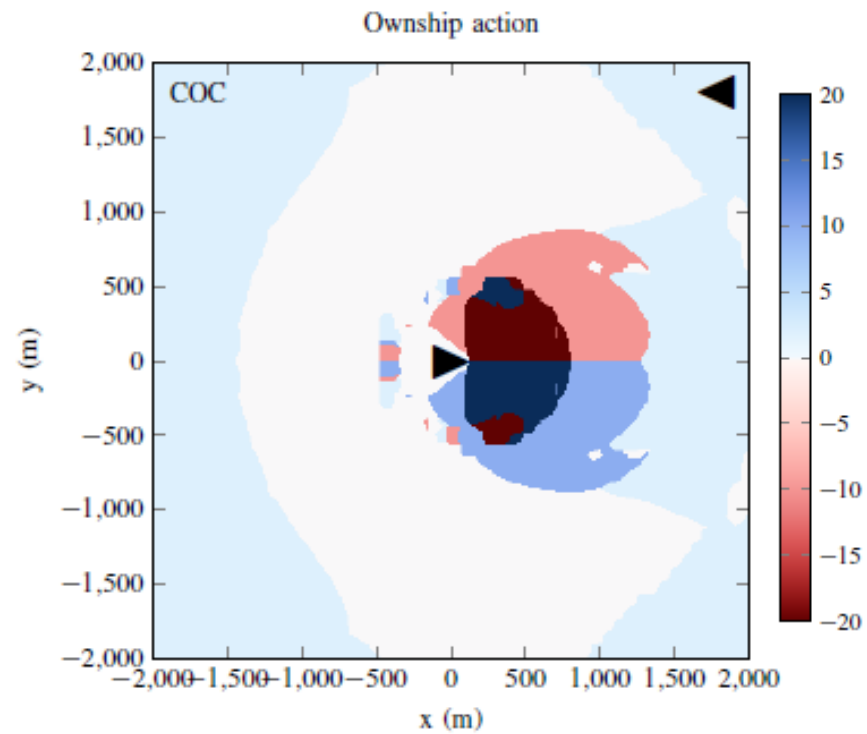
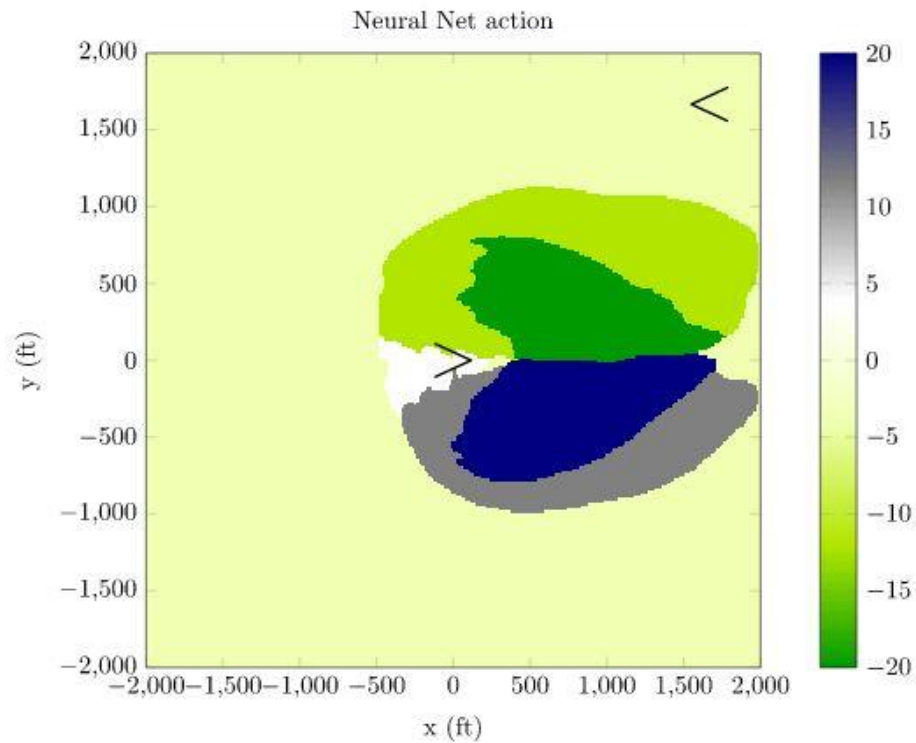
---

## DQN Parameters

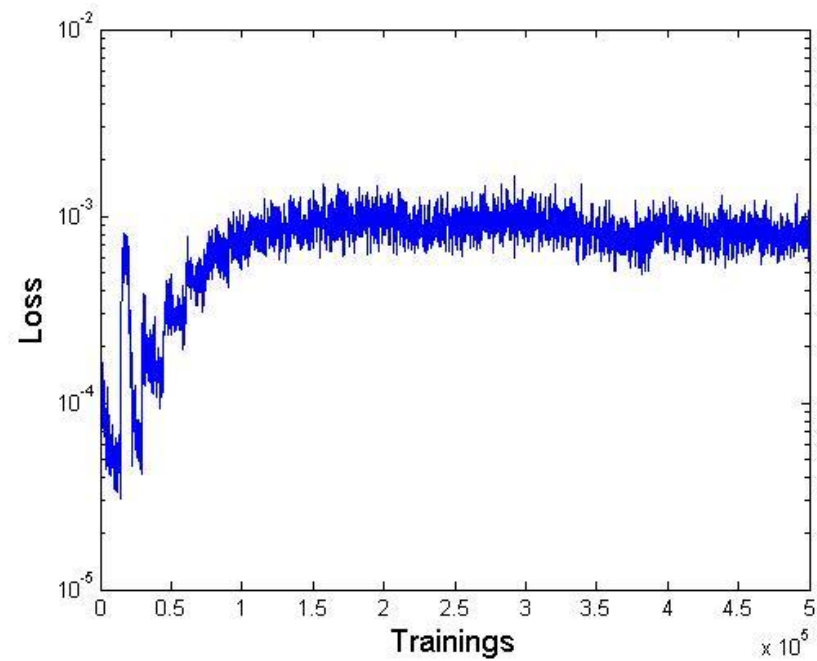
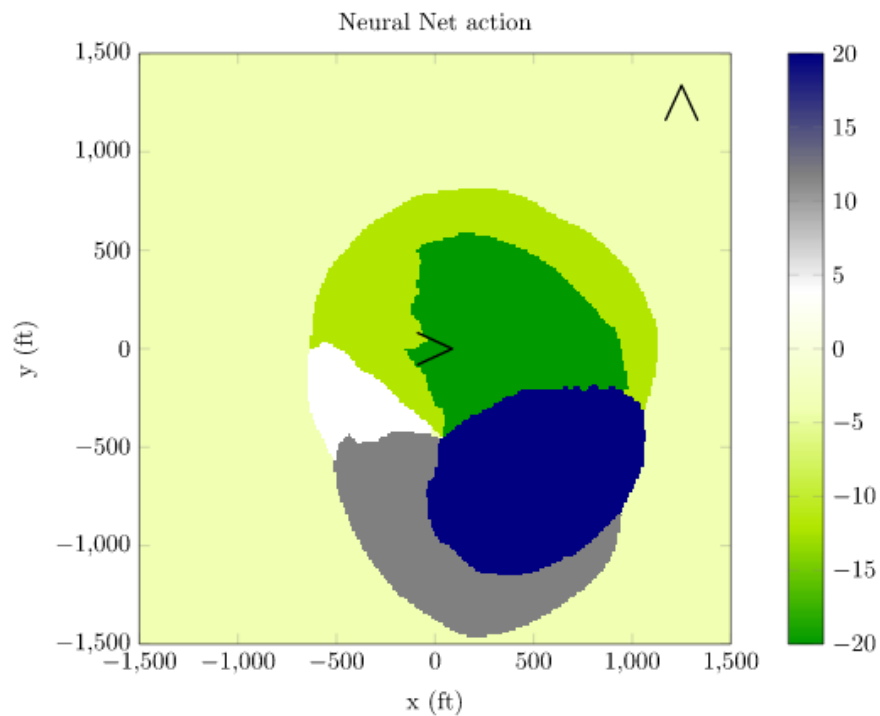
- D\_Initial\_Size = 200,000
- D\_Max\_Size = 1.5M
- Final\_Epsilon = 0.1
- Final\_Exploration = 7.5M
- Train\_Freq = 20
- Batch\_Size = 512
- Target\_Freq = 15000
- Num\_Traj = 16
- Gamma = 0.97
- dt = 5
- Network Library = Keras
- Solver = AdaMax
- Loss = MSE



# Results



# Results



## Future Work

- Prioritized experience replay
- Optimize speed of algorithm
- Use learned parameters with different rewards
- Expand state space
- Simulate policy

# References

- Hao Yi Ong and M. J. Kochenderfer, "Short-term conflict resolution for unmanned aircraft traffic management," *Digital Avionics Systems Conference (DASC), 2015 IEEE/AIAA 34th*, Prague, 2015, pp. 5A4-1-5A4-13.
- Mnih, V. et al. Human-level control through deep reinforcement learning. *Nature* 518, 529–533 (2015).
- Sascha Lange, Thomas Gabel, and Martin Riedmiller. "Batch Reinforcement Learning". M. Wiering, M. Van Otterlo (Eds.). *Reinforcement Learning: State of the Art*. Springer-Verlag (2012). pp. 45-70
- Diederik Kingma and Jimmy Lei Ba. "Adam: A Method for Stochastic Optimization". *International Conference on Learning Representations*. San Diego, 2015.
- M.J. Kochenderfer and J.P. Chryssanthacopoulos, "Robust Airborne Collision Avoidance through Dynamic Programming," MIT Lincoln Laboratory, Project Report ATC-371, 2011.
- Blake Wulfe. Hierarchical Reinforcement Learning Algorithms.  
[https://github.com/wulfebw/hierarchical\\_rl](https://github.com/wulfebw/hierarchical_rl)

Thanks!  
Questions?