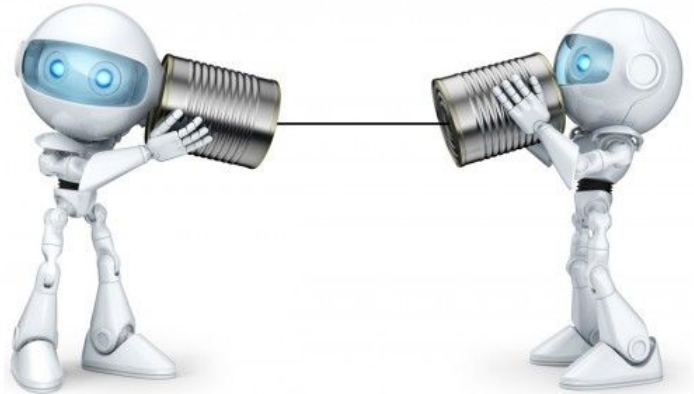


# Training Multiple Intelligent Agents to Communicate

Varun Bhatt

MSc; supervised by Prof. Michael Buro



# Examples of Human Communication



Sharing observations



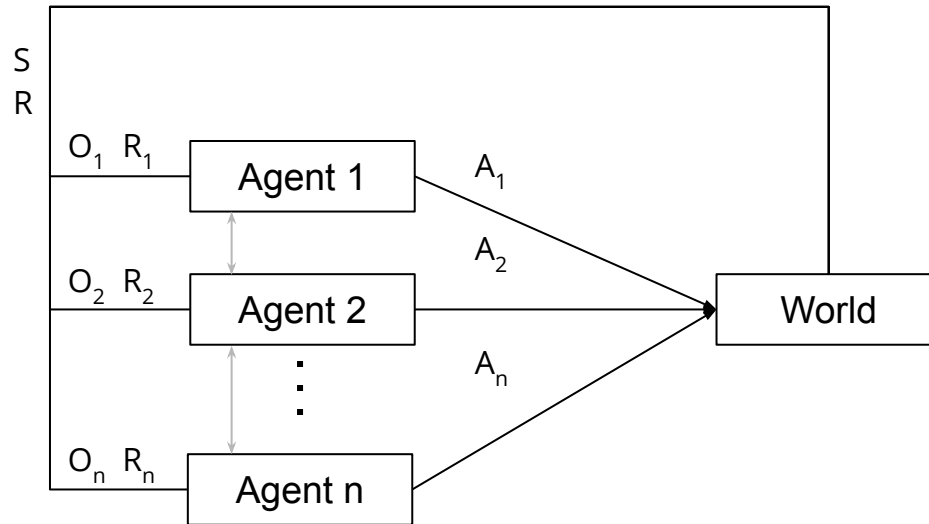
Suggesting action



Sharing action intention

# Formalism

Fully cooperative, partially observable markov game with cheap talk channel  
(dec-POMDP)



# Previous work

Pre-defined communication protocol  
(Tan, 1993)

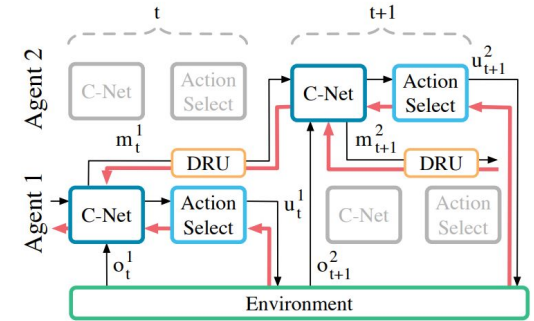


Game theory, equilibrium  
(Claus & Boutilier, 1998)

Agent 2's action

Agent 1's action	11	-30
	-30	7

End to end differentiable communication channel  
(Foerster et. al., 2016)



Tan, M. (1993). Multi-agent reinforcement learning: Independent vs. cooperative agents. In Proceedings of the Tenth International Conference on Machine Learning (pp. 330-337).

Claus, C., & Boutilier, C. (1998). The dynamics of reinforcement learning in cooperative multiagent systems. AAAI/IAAI, 1998, 746-752.

Foerster, J., Assael, I. A., de Freitas, N., & Whiteson, S. (2016). Learning to communicate with deep multi-agent reinforcement learning. In Advances in Neural Information Processing Systems (pp. 2137-2145).

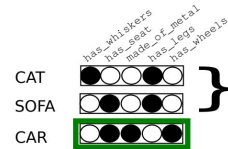
# Previous work

Communication only  
through actions  
(Bard et. al., 2019)



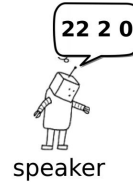
Emergent  
Communication  
(Lazaridou et. al., 2018)

symbolic data

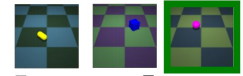


distractors

target

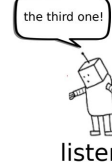


pixel data



distractors

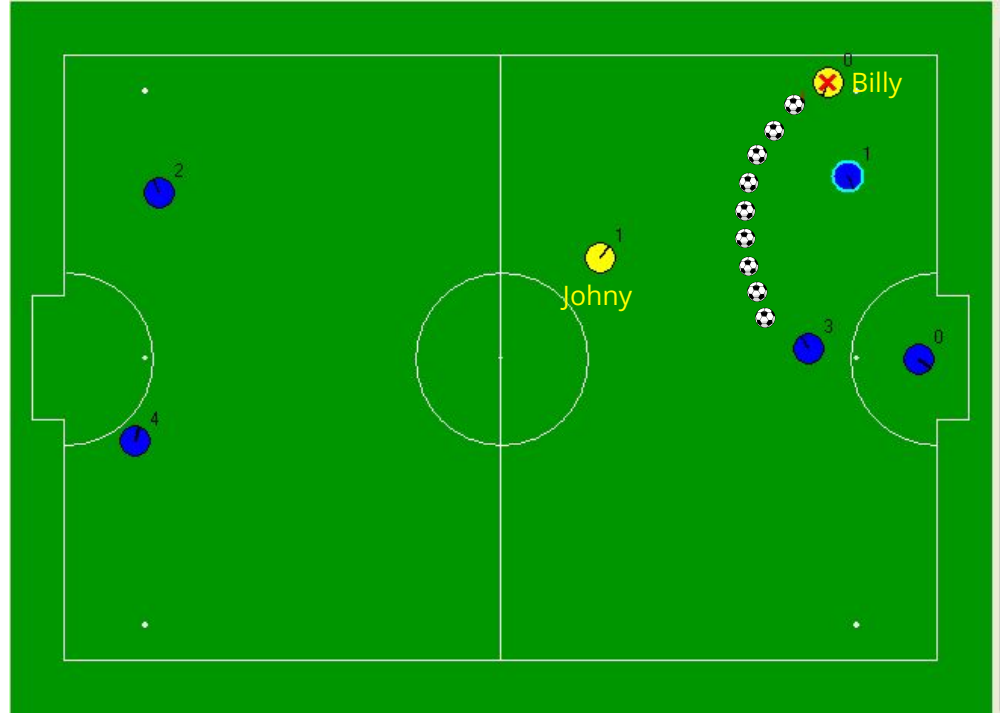
target



# Issues - Credit Assignment

Who's to blame?

Johnny for the call?  
or Billy for the bad pass?



# Issues - Partial Observability

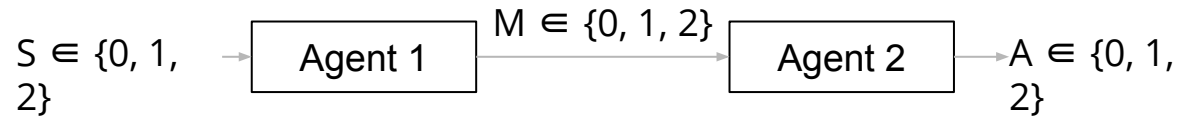


Agent 1

- Can't see  $O_2$
- Doesn't know if  $M_1$  was used or ignored

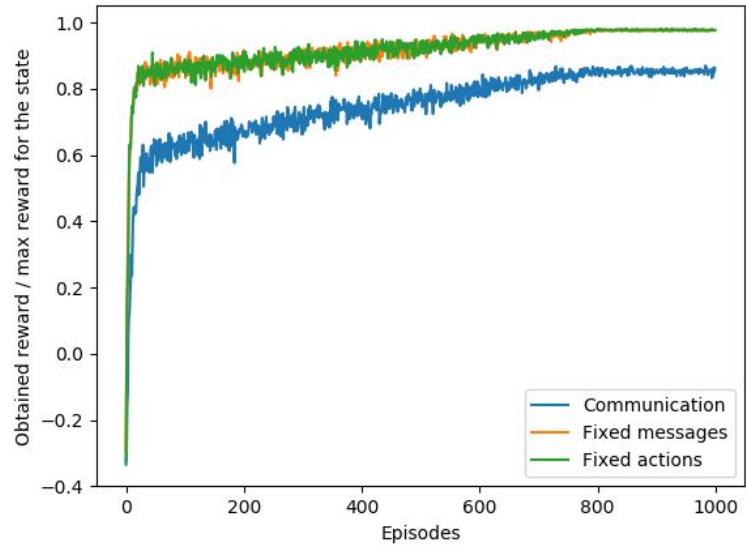
# A Simple Experiment - Modified Climbing Game

	Agent 2's action		
Agent 1's action	11	-30	0
	-30	7	6
	0	0	5

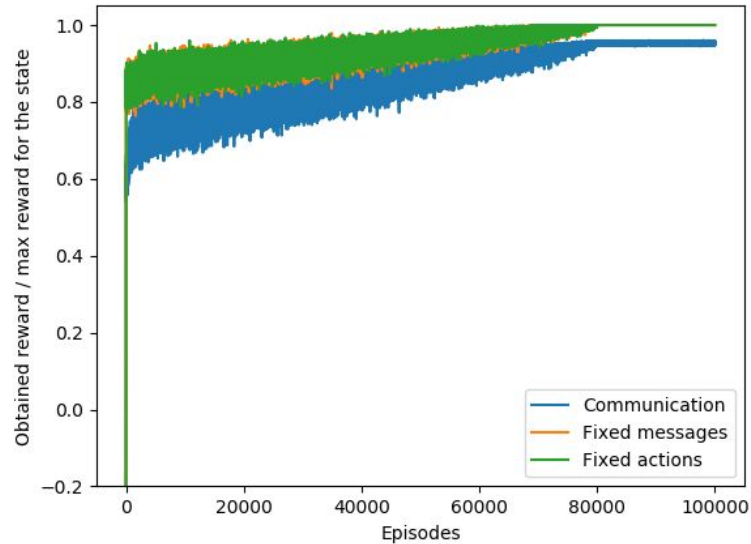




# Results - Reward



# Results - Reward

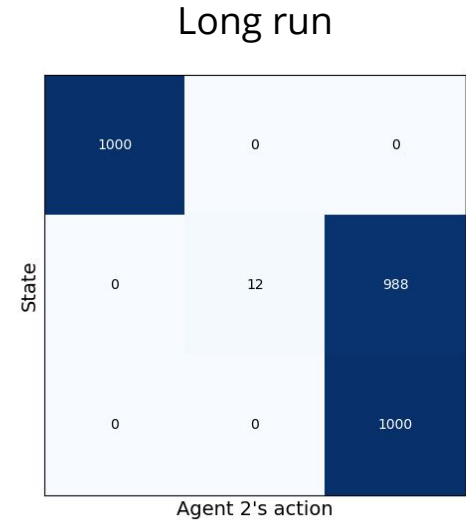
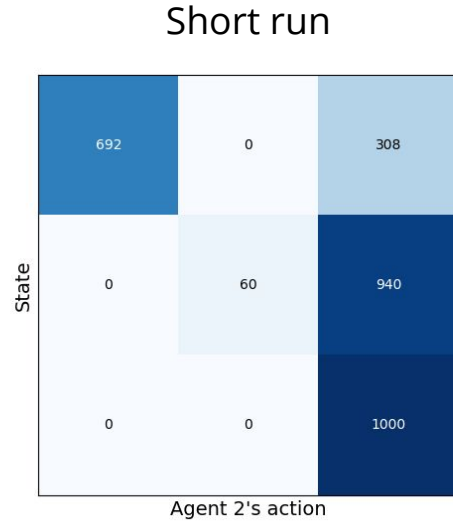


# Results - Convergence points

Agent 2's action

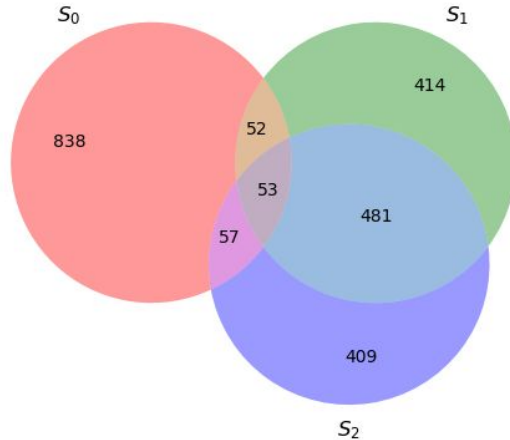
11	-30	0
-30	7	6
0	0	5

State

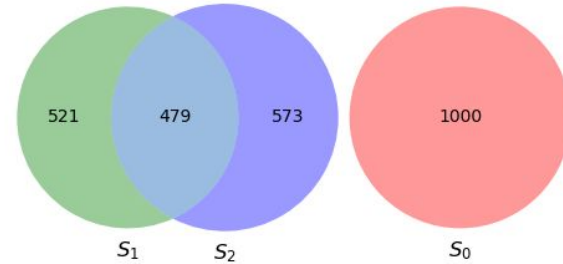


# Results - Message protocol

Short run



Long run

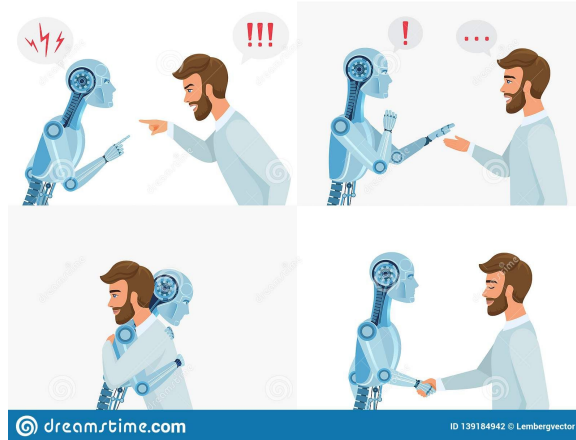


Message is not unique to state

# Ideas

Two timescale optimization/exploration

Two way communication



Questions?