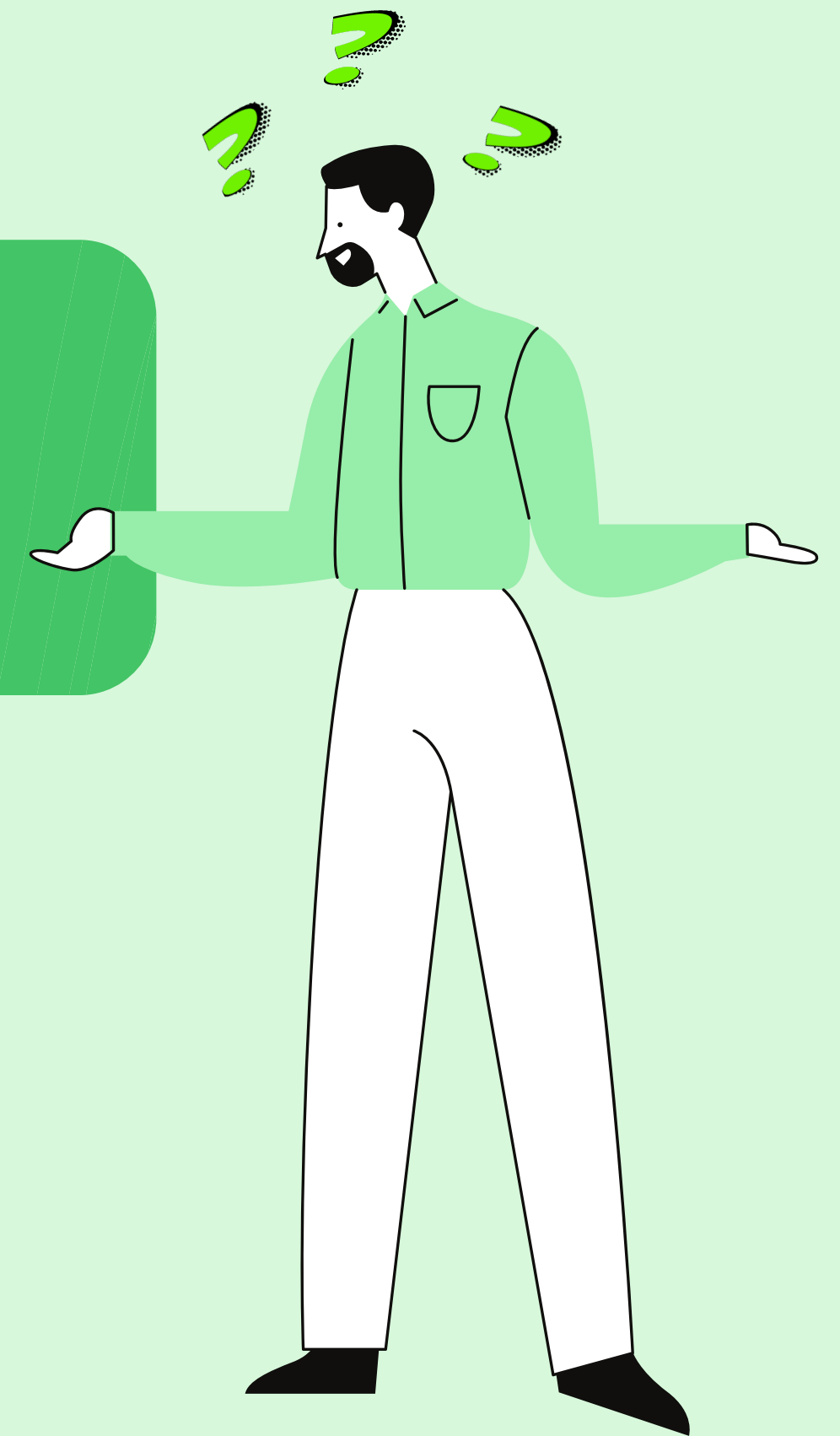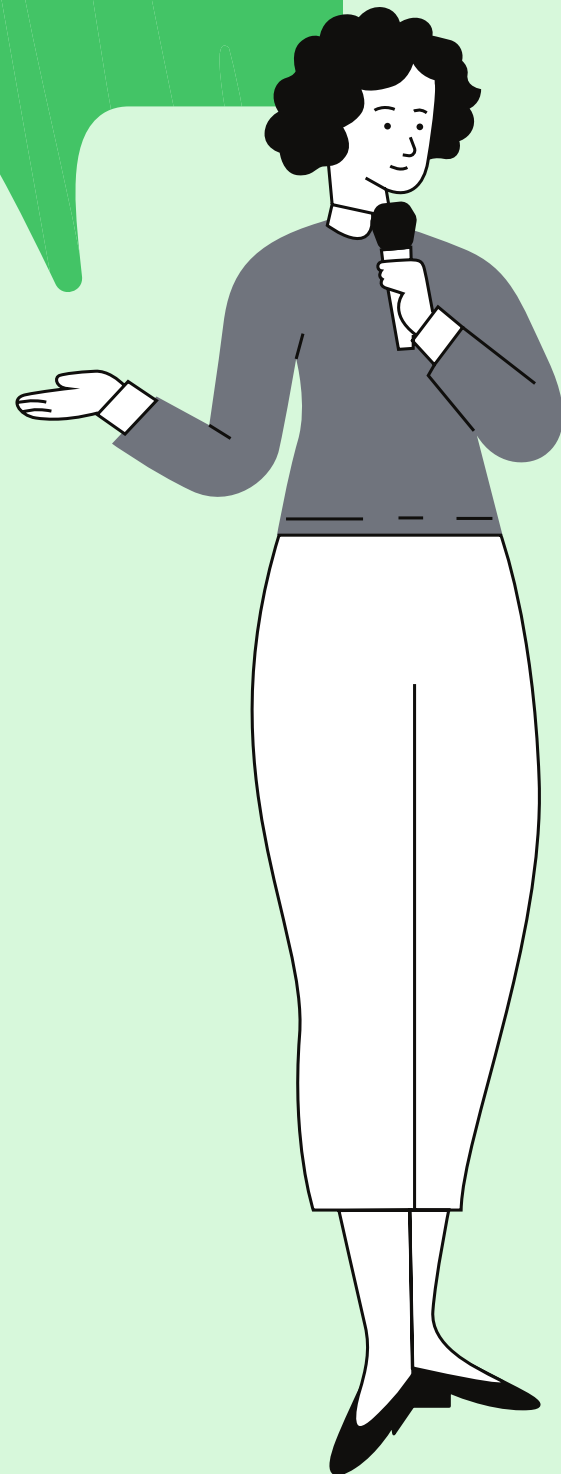# Multi-Agent Systems on Sensor Networks: A Distributed Reinforcement Learning Approach
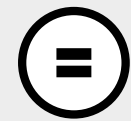
A further breakdown on:
- **Distributed Value Function DRL**
- **Optimistic DRL**

# Distributed Value Function

*Let's start with the general idea behind DVF DRL*

- Nodes communicated by exchanging information about their value functions
- Sum of the individual node value functions

( = )

- Total of future weighted reward sums from all of the nodes
- Nodes have access to the local rewards and can communicate their value functions with their neighbors

## DVF Algorithm

discount factor

learning rate

land in, i.e. $V^i(s_t^i)$ for agent $i$ at time $t$ at each iteration. The update rule at time step $t$ for agent $i$ are given by:

$$Q_{t+1}^i(s_t^i, a_t^i) = (1 - \alpha)Q_t^i(s_t^i, a_t^i) + \alpha\left(r_{t+1}^i(s_{t+1}^i)\right.$$

$$\left. + \gamma \sum_{j \in Neigh(i)} f^i(j)V_t^j(s_{t+1}^j)\right) \tag{2}$$

$$V_{t+1}^i(s_t^i) = \max_{a \in A^i} Q_{t+1}^i(s_t^i, a) \tag{3}$$
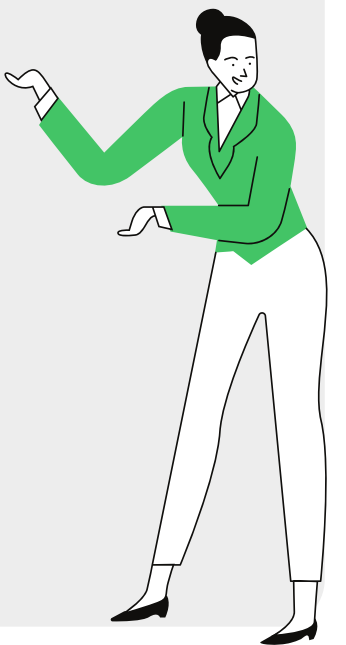
**Reward function /per agent**

$$r^i(s^i) = G^i(s^i) - C^i$$

# Optimistic DRL

*Let's start with the general idea behind Optimistic DRL*

**Generally...**

- All agents are given the same reward function

  - Project **all** information about the whole system as a 'single agent'

- Finds optimal policies in deterministic environments

  - Deterministic is very predictable, your next state is completely dependent on your current state

- Not so great in Stochastic environments

  - Continuous variable –> not exactly predictable

**OptDRL Algorithm**

**Global State**

$$\mathcal{S} = \prod_{i=1}^{m} S^i$$
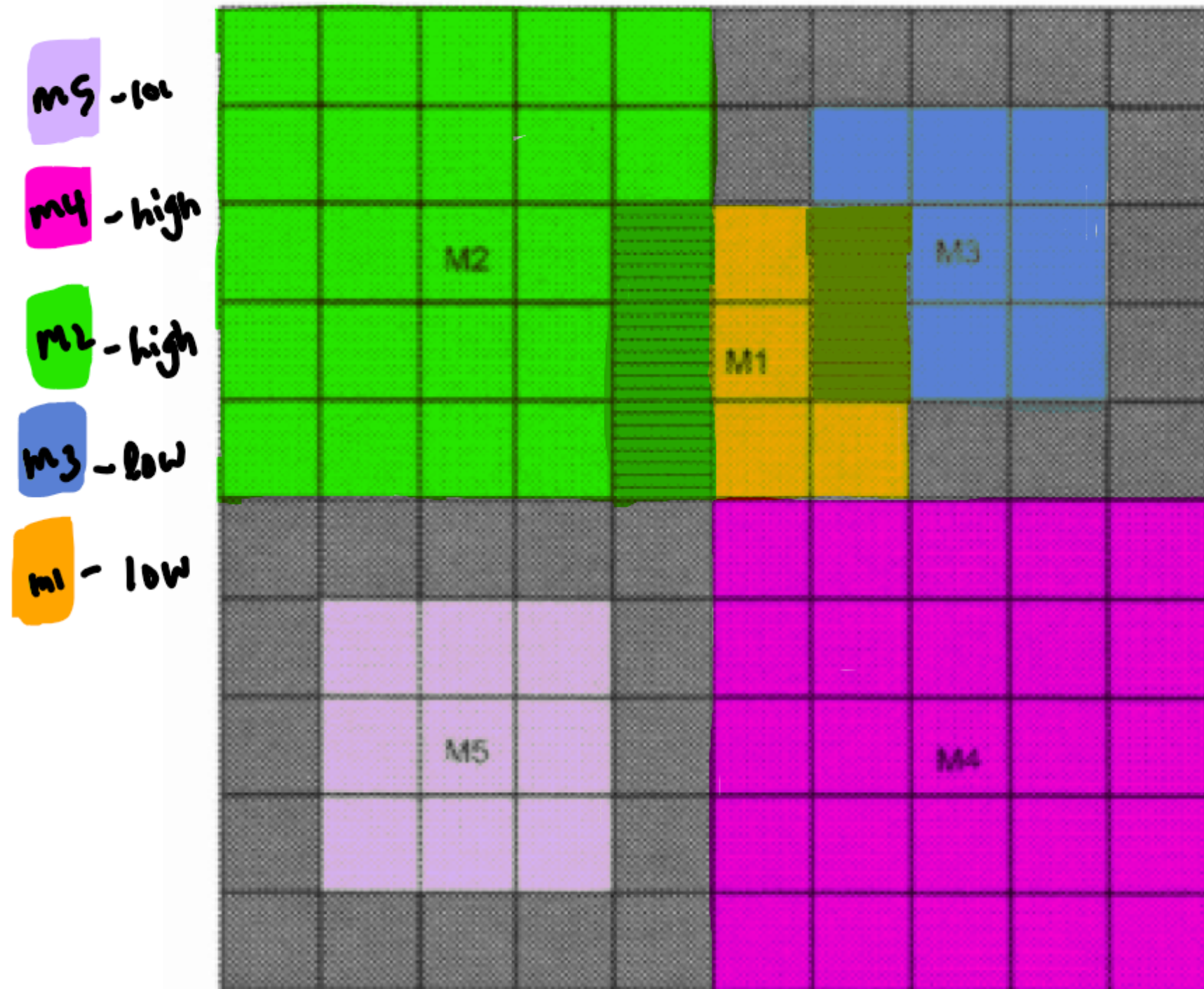
$$q_{t+1}^i(S_t, a_t^i) = \max\{q_t^i(S_t, a_t^i), (1-\alpha)q_t^i(S_t, a_t^i) + \alpha(r_{t+1}^i(S_{t+1}) + \gamma \max_{a \in A^i} q_t^i(S_{t+1}, a))\} \quad (5)$$

$$\Pi_{t+1}^i(S_t) = a_t^i \text{ iff } \max_{a \in A^i} q_t^i(S_t, a) \neq \max_{a \in A^i} q_{t+1}^i(S_t, a) \quad (6)$$

**Reward function**
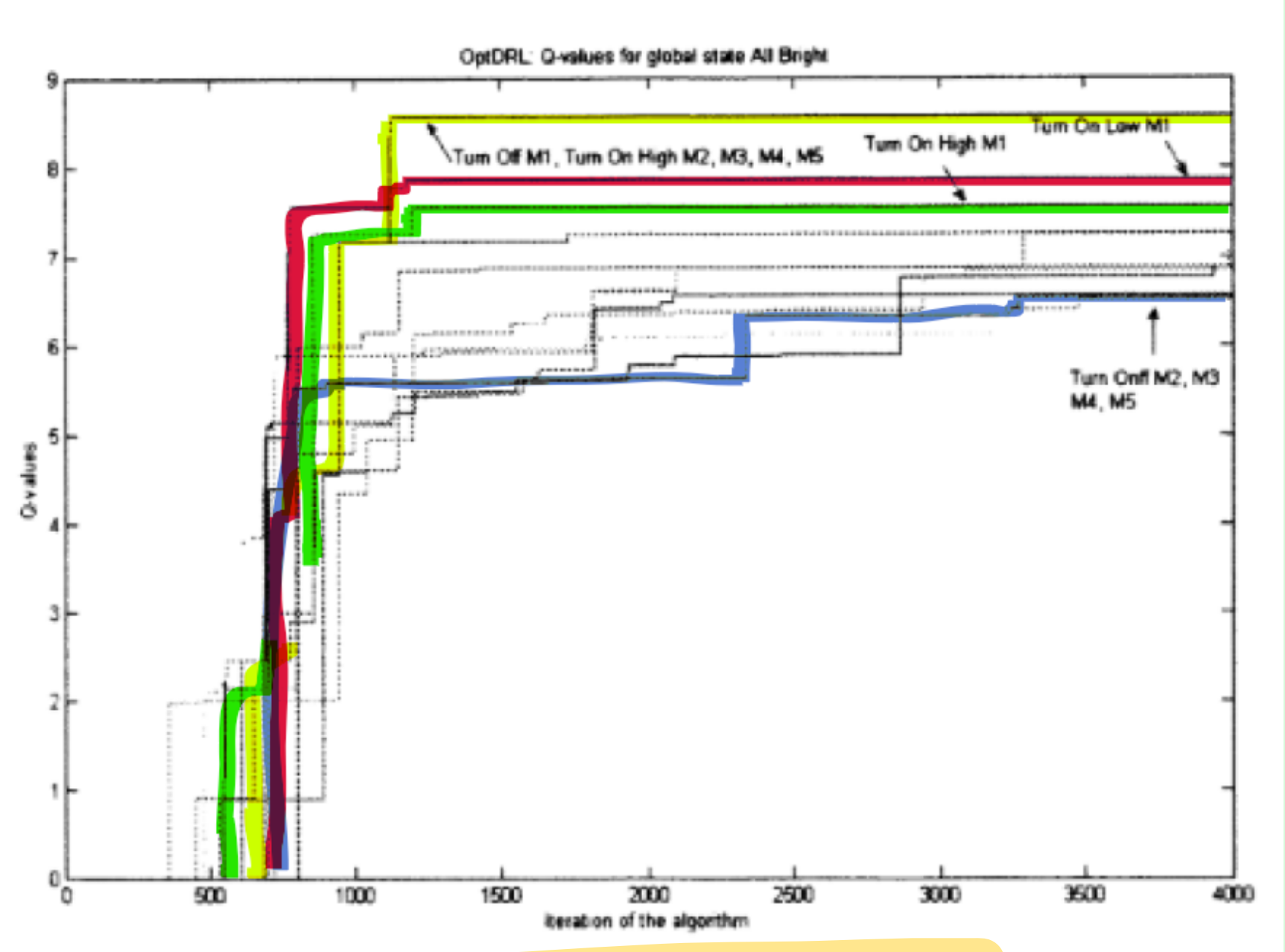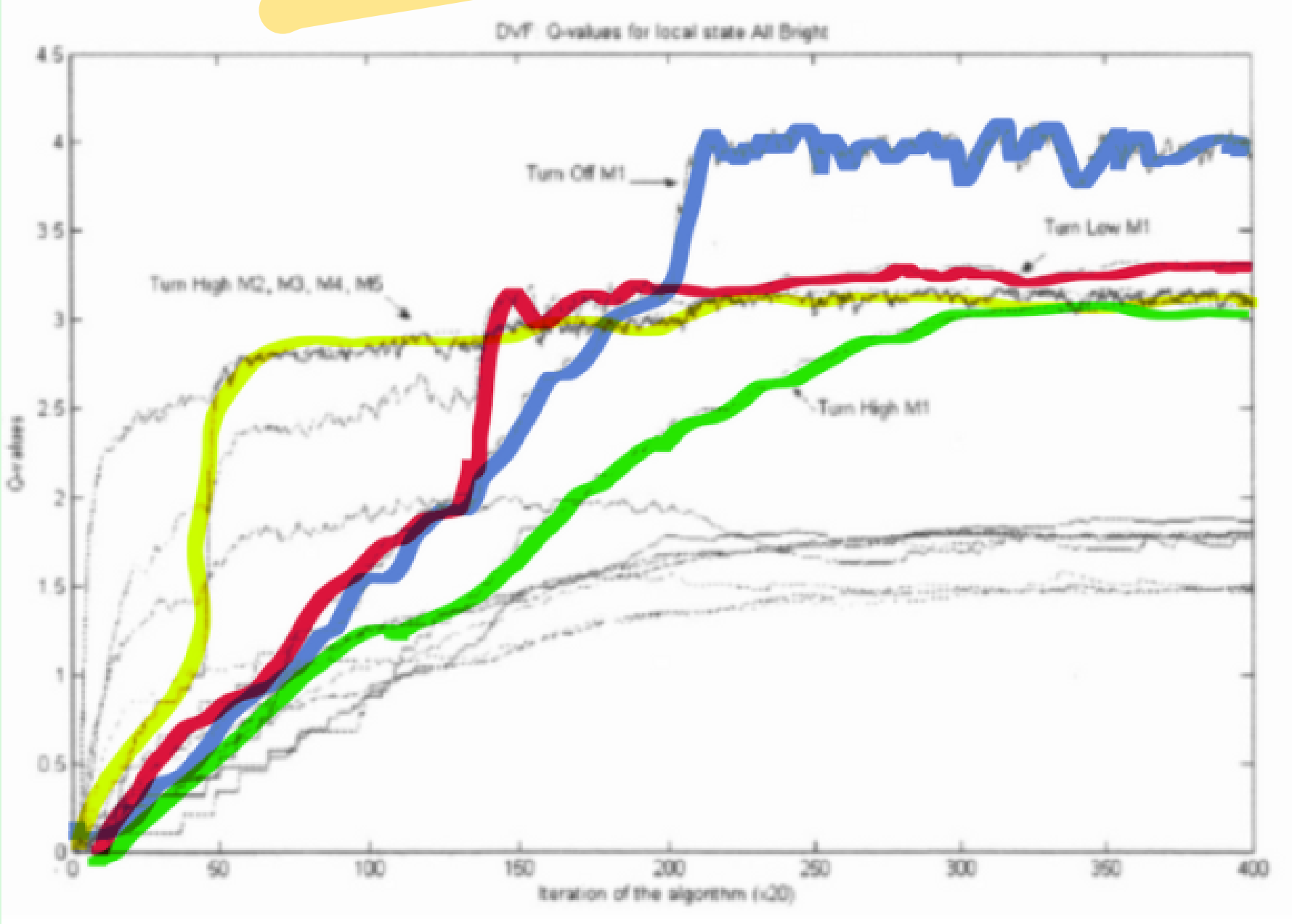
$$Glob\_rew(S) = G(S) - C(S)$$

# Simulation



**Setting**

- 10x10 Grid
- 5 Agents
  - Off – no lights
  - low – 9 squares
  - high – 25 squares
- Goal: light up ALL 100 squares as efficiently as possible

# Simulation Results

DVF DRL



DVF: Q-values for local state All Bright

Turn Off M1

Turn Low M1

Turn High M2, M3, M4, M5

Turn High M1

Iteration of the algorithm (x20)

OPT DRL



OptDRL: Q-values for global state All Bright

Turn Off M1, Turn On High M2, M3, M4, M5

Turn On High M1

Turn On Low M1

Turn Onff M2, M3 M4, M5

Iteration of the algorithm

# Simulation Results

TABLE 1: ENERGY CONSUMPTION (J) AND APPLICATION-LEVEL PE
FORMANCE OF THE MAS WITH 5 MOTES DURING THE FIRST 4000 ITER
TIONS

|  | Ind Learners | DVF | OptDRL |
|---|---|---|---|
| Communication | 0 | 1648.4 | 1681.0 |
| Computation | 954.9 | 971.7 | 1187.8 |
| Lights LOW | 2404 | 3014 | 1393 |
| Lights HIGH | 12721 | 12192 | 12123 |
| Cells Bright | 318323 | 317674 | 321766 |
| Cumulative Reward | 17429 | 17466 | 19078 |

TABLE 2: MEMORY REQUIREMENTS OF THE ALGORITHMS

|  | Expression | Actual values |
|---|---|---|
| IndLearners | $|s^i| \times |A^i|$ | $2^{25} \times 3$ |
| DVF | $|s^i| \times |A^i| + |s^i|$ | $2^{25} \times 4$ |
| OptDRL | $|S| \times |A^i| + |S|$ | $2^{100} \times 4$ |

# Integrating it with Current Project

- The use of multi agent to converge to an optimal policy does seem promising
  - Could use multiple agents to help converge to the reach the lowest communication cost
- A potential issue is looking for real live coding examples
  - different terms
  - perhaps only present in theoretical studies
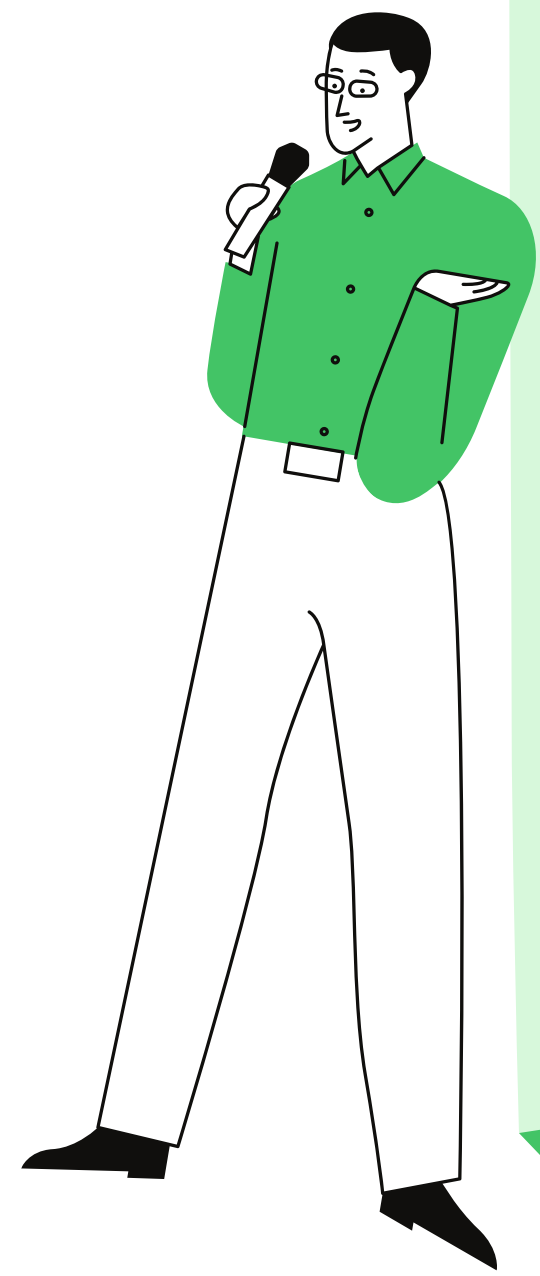  - requires further digging for resources

# Works Cited

Chen-Khong Tham and J. -. Renaud, "Multi-Agent Systems on Sensor Networks: A Distributed Reinforcement Learning Approach," 2005 International Conference on Intelligent Sensors, Sensor Networks and Information Processing, 2005, pp. 423-429, doi: 10.1109/ISSNIP.2005.1595616.