

Multi-agent Determinantal Q-Learning

Yaodong Yang^{*1,2}, Ying Wen^{*1,2}, Liheng Chen³, Jun Wang², Kun Shao¹, David Mguni¹, Weinan Zhang³

¹*Huawei Technology R&D UK*

²*University College London*

³*Shanghai Jiao Tong University*



06/2020

Multi-Agent Reinforcement Learning for Real-World Applications

- MARL answers the question of how each agent should behave optimally when interacting with the others.



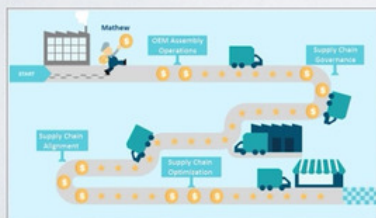
Gaming AI



Robotics



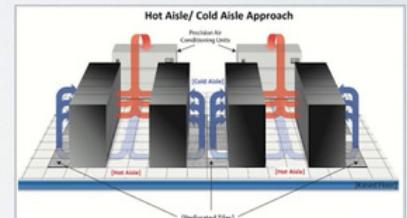
Autonomous driving



Supply chain optimisation



Network control



Data centre cooling

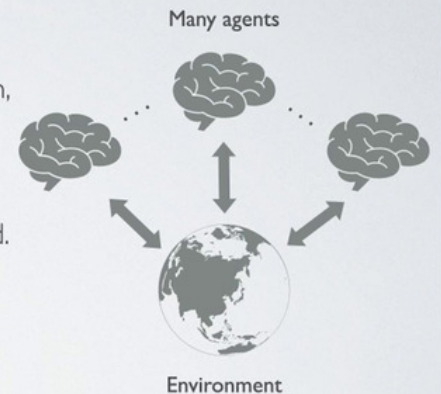
Problem Formulation: Multi-agent Reinforcement Learning

- Modelled by a Stochastic Game $(\mathcal{S}, \mathcal{A}^{\{1,\dots,n\}}, \mathcal{R}^{\{1,\dots,n\}}, \mathcal{T}, \mathcal{P}_0, \gamma)$
 - \mathcal{S} denotes the state space,
 - \mathcal{A} is the joint-action space $\mathcal{A}^1 \times \dots \times \mathcal{A}^n$,
 - $\mathcal{R}^i = \mathcal{R}^i(s, a^i, a^{-i})$ is the reward function for the i-th agent,
 - $\mathcal{T} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0,1]$ is the transition function based on the joint action,
 - \mathcal{P}_0 is the distribution of the initial state, γ is a discount factor.
 - **Special case:** $n = 1 \rightarrow$ single-agent MDP, $|\mathcal{S}| = 1 \rightarrow$ normal-form game
 - **Dec-POMDP:** state is only partially observed, agents share the same reward.

- Each agent tries to maximise its expected long-term reward:

$$V_{i,\pi}(s) = \sum_{t=0}^{\infty} \gamma^t \mathbf{E}_{\pi, \mathcal{P}} \{ R_{i,t} | s_0 = s, \pi \}, \pi = [\pi_1, \dots, \pi_N]$$

$$Q_{j,\pi}(s, a) = R_j(s, a) + \gamma \mathbf{E}_{s' \sim p} [v_{j,\pi}(s')]$$



Solution to Multi-Agent RL

- Value-based method:

- The sense of optimality changes, now it depends on other agents !

$$\pi_{i,t}(s, \cdot) = \text{solve}_i \{ Q_{\cdot,t}(s_t, \cdot) \}$$

$$Q_{i,t+1}(s_t, \pi_t) = Q_{i,t}(s_t, \pi_t) + \alpha [R_{i,t+1} + \gamma \cdot \text{eval}_i \{ Q_{\cdot,t}(s_{t+1}, \cdot) \} - Q_{i,t}(s_t, \pi_t)]$$

- ♦ Fully-cooperative game: agents share the same reward function

$$\text{solve}_i \{ Q_{\cdot,t}(s_t, \cdot) \} = \arg \max_{a_i} \left(\max_{a^{-i}} Q_{i,t}(s_t, a^i, a^{-i}) \right)$$

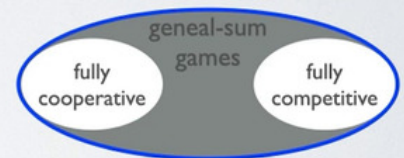
$$\text{eval}_i \{ Q_{\cdot,t}(s_{t+1}, \cdot) \} = \max_a Q_{i,t}(s_{t+1}, a)$$

- ♦ Fully-competitive game: sum of agents' reward is zero

$$\text{solve}_i \{ Q_{\cdot,t}(s_t, \cdot) \} = \arg \max_{a_i} \left(\min_{a^{-i}} Q_{i,t}(s_t, a^i, a^{-i}) \right)$$

$$\text{eval}_i \{ Q_{\cdot,t}(s_{t+1}, \cdot) \} = \max_{a_i} \min_{a^{-i}} Q_{i,t}(s_{t+1}, a^i, a^{-i})$$

- In this work, we focus on cooperative MARL with partial observability.



Cooperative Multi-agent Reinforcement Learning

- Decentralising the joint Q-function is an essential approach to scale MARL.

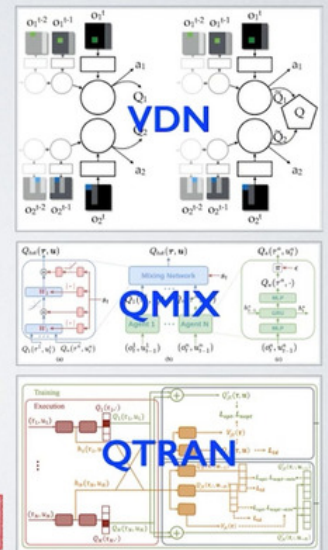
- A task is decentralisable if

$$\arg \max_a Q^\pi(s, a) = \begin{bmatrix} \arg \max_{a_1} Q_1(s_1, a_1) \\ \vdots \\ \arg \max_{a_N} Q_N(s_N, a_N) \end{bmatrix}$$

- VDN [Sunehag 2017]: $Q^\pi(s, a) := \sum_{i=1}^N Q_i(s_i, a_i)$

- QMIX [Rashid 2018]: $\frac{\partial Q^\pi(s, a)}{\partial Q_i(s_i, a_i)} \geq 0, \forall i \in \mathcal{N}$

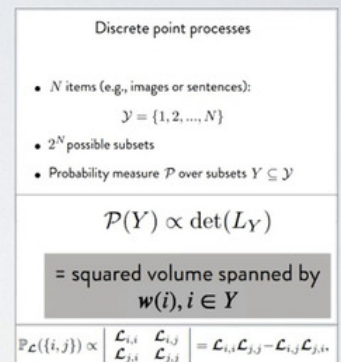
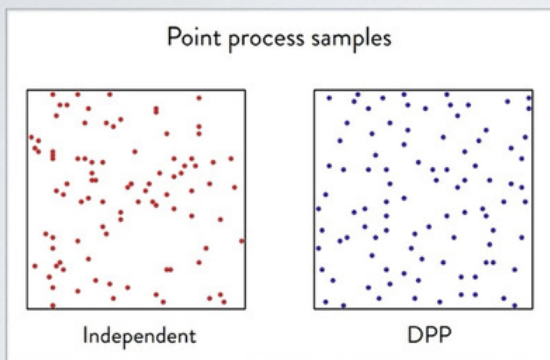
- QTRAN [Kyunghwan 2019]: $\sum_{i=1}^N Q_i(s_i, a_i) - Q^\pi(s, a) + V_\pi(a) = 0$ if $a = a^*$
 ≥ 0 if $a \neq a^*$



Can we remove these special structural constraints on $Q^\pi(s, a)$?

Determinantal Point Process: A Probabilistic Framework for Diversity

- Intuition: if the agents learn diverse behaviours, then $Q^\pi(s, a)$ can be naturally decentralised.
 - Determinantal Point Process (DPP) [Alex Kulesza 2013] : a point process parameterised by a distance kernel.



$$\mathbb{P}_{\mathcal{L}}(Y = Y) \propto \det(\mathcal{L}_Y) = \text{Vol}^2(\{w_i\}_{i \in Y})$$

Q-DPP: A New Function Approximator for Cooperative MARL

- **Q-Determinantal Point Process (Q-DPP):** a new function approximator for the joint Q-function.
 - ♦ The DPP kernel \mathcal{L} can represent the joint Q-function $Q^\pi(o, a)$.

Definition 3 (Q-DPP). Given a ground set \mathcal{Y} of size M that includes N agents' all possible observation-action pairs $\mathcal{Y} = \{(o_1^1, a_1^1), \dots, (o_N^1, a_N^1)\}$, we partition \mathcal{Y} into N disjoint parts, i.e., $\mathcal{Y} = \bigcup_{i=1}^N \mathcal{Y}_i$ and $\sum_{i=1}^N |\mathcal{Y}_i| = M = N|\mathcal{O}||\mathcal{A}|$, where each partition represents each individual agent's all possible observation-action pairs. At every time-step, given agents' observations, $\mathbf{o} = (o_i)_{i \in \mathcal{N}}$, we define $\mathcal{C}(\mathbf{o}) \subseteq \mathcal{Y}$ to be a set of valid subsets including only observation-action pairs that agents are allowed to take,

$$\mathcal{C}(\mathbf{o}) := \{Y \subseteq \mathcal{Y} : |Y \cap \mathcal{Y}_i(o_i)| = 1, \forall i \in \{1, \dots, N\}\},$$

with $|\mathcal{C}(\mathbf{o})| = |\mathcal{A}|^N$, and $\mathcal{Y}_i(o_i)$ of size $|\mathcal{A}|$ denotes the set of pairs in partition \mathcal{Y}_i with only o_i as the observation,

$$\mathcal{Y}_i(o_i) = \{(o_i, a_1^i), \dots, (o_i, a_{|\mathcal{A}|}^i)\}.$$

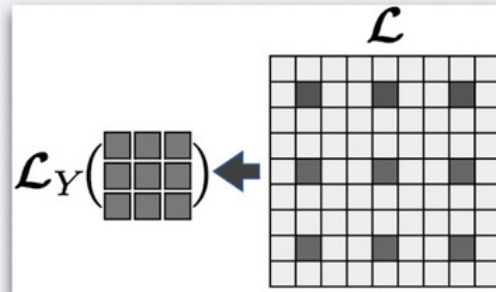
Q-DPP, denoted by $\tilde{\mathbb{P}}$, defines a probability measure over the valid subsets $Y \in \mathcal{C}(\mathbf{o}) \subseteq \mathcal{Y}$. Let \mathbf{Y} be a random subset drawn according to $\tilde{\mathbb{P}}$, its probability distribution is defined:

$$\tilde{\mathbb{P}}_{\mathcal{L}}(\mathbf{Y} = Y | \mathbf{Y} \in \mathcal{C}(\mathbf{o})) := \frac{\det(\mathcal{L}_Y)}{\sum_{Y' \in \mathcal{C}(\mathbf{o})} \det(\mathcal{L}_{Y'})}. \quad (4)$$

In addition, given a valid sample $Y \in \mathcal{C}(\mathbf{o})$, we define an identifying function $\mathcal{I} : Y \rightarrow \mathcal{N}$ that specifies the agent number for each valid pair in Y , and an index function $\mathcal{J} : \mathcal{Y} \rightarrow \{1, \dots, M\}$ that specifies the cardinality of each item in Y in the ground set \mathcal{Y} .

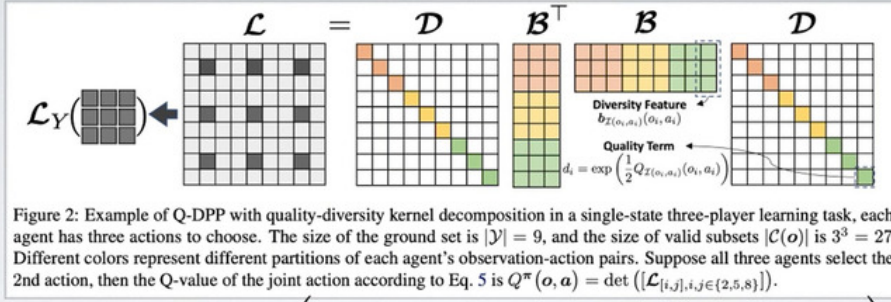
$$Q^\pi(o, a) := \log \det \left(\mathcal{L}_{Y=\{(o_1^1, a_1^1), \dots, (o_N^1, a_N^1)\} \in \mathcal{C}(\mathbf{o}^1)} \right)$$

Example: single-state three-player three-action task, the probability of all agents select their second action is represented by the following sampling process of DPP.



Q-DPP: A New Function Approximator for Cooperative MARL

- Q-Determinantal Point Process (Q-DPP): a new function approximator for the joint Q-function.
 - The DPP kernel \mathcal{L} can represent both quality (maximising reward) & diversity (different behaviours).



$$\begin{aligned}
 Q^\pi(o, a) &:= \log \det \left(\mathcal{L}_{Y=\{(o_1^t, a_1), \dots, (o_N^t, a_N)\} \in \mathcal{C}(o^t)} = \mathcal{W}_Y \mathcal{W}_Y^\top \right) \\
 &= \log \left(\text{tr}(\mathcal{D}_Y^\top \mathcal{D}_Y) \det(\mathcal{B}_Y^\top \mathcal{B}_Y) \right) \quad \text{quality-diversity decomposition} \\
 &= \sum_{i=1}^N \underbrace{Q_{\mathcal{J}(o_i^t, a_i)}(o_i, a_i)}_{\text{quality}} + \log \det(\mathcal{B}_Y^\top \mathcal{B}_Y) \quad \text{diversity}
 \end{aligned}$$

Sampling from Q-DPP through Orthogonal Projection

- Q-DPP Sampler: each agent's decision making is equivalent to sampling from the Q-DPP.

- ♦ Kernel \mathcal{L} can be learned by Q-learning, but sampling from Q-DPP is challenging since $|\mathcal{C}(\theta)| = |\mathcal{A}|^N$.

$$\tilde{\mathbb{P}}_{\mathcal{L}}(Y = Y | Y \in \mathcal{C}(\theta)) := \frac{\det(\mathcal{W}_Y \mathcal{W}_Y^\top)}{\sum_{Y' \in \mathcal{C}(\theta)} \det(\mathcal{W}_{Y'} \mathcal{W}_{Y'}^\top)}$$

- ♦ Best DPP sampler needs $\mathcal{O}(|\mathcal{A}|^{N^2})$, we instead propose a biased sampler and then control the bias.
 - Sample each partition sequentially, then the size of ground set is $|\mathcal{C}(\theta)| = |\mathcal{A}|^N \rightarrow |\mathcal{A}|N$.
 - The validity comes from the fact that **Gram-Schmidt process** preserve the volume.

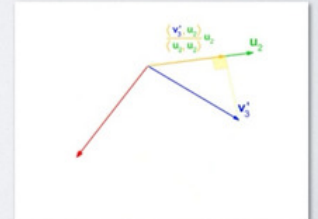
$$\det(\mathcal{W}_Y \mathcal{W}_Y^\top) = \prod_{i=1}^M \|\mathbf{U}_{\mathcal{U}_i}(\mathbf{w}_i)\|^2, \mathcal{U}_i = \text{span}\{\mathbf{w}_1, \dots, \mathbf{w}_{i-1}\}$$

- Sampling is also made easy if the items in partitions are orthogonalised.

$$\mathbb{P}(Y) \propto \prod_{i \in Y} \|\mathbf{w}_i\|^2 = \det(\mathcal{W}_Y \mathcal{W}_Y^\top) \propto \det(\mathcal{L}_Y)$$

this is what we want!

- Remaining issue: how to control of the bias in the denominator.



Gram-Schmidt Process [wiki]

Controlling the Bias of the Q-DPP Sampler

- Q-DPP sampler, though is linear-time, is biased in that samples have larger probability.

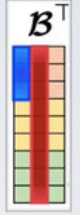
```

1: DEF Orthogonalizing-Sampler ( $\mathcal{Y}, \mathcal{D}, \mathcal{B}, o$ ):
2: Init:  $\mathbf{b}_j \leftarrow \mathcal{B}_{[:,j]}, Y \leftarrow \emptyset, B \leftarrow \emptyset, J \leftarrow \emptyset$ .
3: for each partition  $\mathcal{Y}_i$  do
4:   Define  $\forall (o, a) \in \mathcal{Y}_i(o_i)$ 
      $q(o, a) := \|\mathbf{b}_{\mathcal{J}(o,a)}\|^2 \exp(\mathcal{D}_{\mathcal{J}(o,a), \mathcal{J}(o,a)})$ .
5:   Sample  $(\tilde{o}_i, \tilde{a}_i) \in \mathcal{Y}_i(o_i)$  from the distribution:
     
$$\left\{ \frac{q(o, a)}{\sum_{(\hat{o}, \hat{a}) \in \mathcal{Y}_i(o_i)} q(\hat{o}, \hat{a})} \right\}_{(o,a) \in \mathcal{Y}_i(o_i)}$$

6:   Let  $Y \leftarrow Y \cup (\tilde{o}_i, \tilde{a}_i), B \leftarrow B \cup \mathbf{b}_{\mathcal{J}(\tilde{o}_i, \tilde{a}_i)},$ 
      $J \leftarrow J \cup \mathcal{J}(\tilde{o}_i, \tilde{a}_i)$ .
7:   // Gram-Schmidt orthogonalization
8:   Set  $\mathbf{b}_j = \Pi_{\text{span}\{B\}}(\mathbf{b}_j), \forall j \in \{1, \dots, M\} - J$ 
9: end for
10: Return:  $Y$ .
    
```

- Q-DPP defines the denominator $\sum_{Y \in \mathcal{C}(o)} \det(\mathcal{W}_Y \mathcal{W}_Y^\top)$.
- Inspired by P-DPP [Celis 2018], we can control the level of violation by bounding the singular value of each partition, such that, $\mathbb{P}(Y) \leq 1/\delta^N \cdot \tilde{\mathbb{P}}_{\mathcal{G}}(Y = Y)$.

Assumption 1 (Singular-Value Constraint on Partitions).
 For a Q-DPP defined in Definition 1, which is parameterized by $\mathcal{D} \in \mathbb{R}^{M \times M}, \mathcal{B} \in \mathbb{R}^{P \times M}$ and $\mathcal{W} := \mathcal{D}\mathcal{B}^\top$, let $\sigma_1 \geq \dots \geq \sigma_P$ represent the singular values of \mathcal{W} , and $\hat{\sigma}_{i,1} \geq \dots \geq \hat{\sigma}_{i,P}$ denote the singular values of $\mathcal{W}_{\mathcal{Y}_i}$ that is the submatrix of \mathcal{W} with the rows and columns corresponding to the i -th partition \mathcal{Y}_i , we assume $\forall j \in \{1, \dots, P\}, \exists \delta \in (0, 1]$, s.t., $\min_{i \in \{1, \dots, N\}} \hat{\sigma}_{i,j}^2 / \delta \geq \sigma_j^2$ holds.



- Total complexity $\mathcal{O}(N^2 | \mathcal{A} | P)$ with input $\mathcal{O}(N | \mathcal{A} | P)$.
- To maintain such assumption, we adopt an additional loss function of $\max(0, \sigma_i^2 - \hat{\sigma}_{i,i}^2 / \delta)$ on top of the Q-update.

Learning the Q-DPP Kernel through CTDE Method

- Note that $Q^\pi(o, a)$ is still learned via normal Q-learning update:

```

12: DEF Determinantal-Q-Learning ( $\theta = [\theta_{\mathcal{D}}, \theta_{\mathcal{B}}], \mathcal{Y}$ ):
13: Init:  $\theta^- \leftarrow \theta, D \leftarrow \emptyset$ .
14: for each time-step do
15:   Collect observations  $o = [o_1, \dots, o_N]$  for all agents.
16:    $a = \text{Orthogonalizing-Sampler}(\mathcal{Y}, \theta_{\mathcal{D}}, \theta_{\mathcal{B}}, o)$ .
17:   Execute  $a$ , store the transition  $\langle o, a, \mathcal{R}, o' \rangle$  in  $D$ .
18:   Sample a mini-batch of  $\{\langle o, a, \mathcal{R}, o' \rangle\}_{j=1}^E$  from  $D$ .
19:   Compute for each transition in the mini-batch
        $\max_{a'} Q(o', a'; \theta^-)$ 
        $= \log \det (\mathcal{L}_{Y=\{(o'_1, a_1^*), \dots, (o'_N, a_N^*)\}})$ 
       where // off-policy decentralized execution
        $a_i^* = \arg \max_{a_i \in \mathcal{A}_i} \left[ \|\theta_{\mathcal{D}}^-(o'_i, a_i)\|^2 \right.$ 
       decentralized
       execution  $\left. \cdot \exp(\theta_{\mathcal{D}}^-(o'_i, a_i), \mathcal{J}(o'_i, a_i)) \right]$ .
20: // centralized training
21:   Update  $\theta$  by minimizing  $\mathcal{L}(\theta)$  defined in Eq. 2.
22:   Update target  $\theta^- = \theta$  periodically.
23: end for
24: Return:  $\theta_{\mathcal{D}}, \theta_{\mathcal{B}}$ .

```

$$\mathcal{L}(\theta) = \sum_{j=1}^E \left\| \mathcal{R} + \gamma \max_{a'} Q(o', a'; \theta^-) - Q^\pi(o, a; \theta) \right\|^2$$

- Q-DPP is a CTDE method. During training, the joint Q-function is updated. During execution, only each agent's own quality and diversity terms are needed.
- The diversity guarantee ensures the Q-DPP can meet the decentralizability assumption after training.
- The orthogonalizing sampling procedures help agents explore different state-space during training thus boosting the efficiency of exploration.
- Note the Q-DPP can only deal with discrete state-action case. The deep Q-DPP is on its way.

Q-DPP Performance

- Performance: the new SOTA algorithm with no neural network needed!

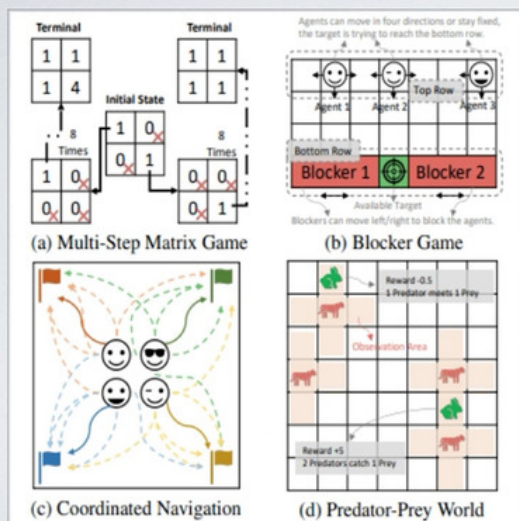
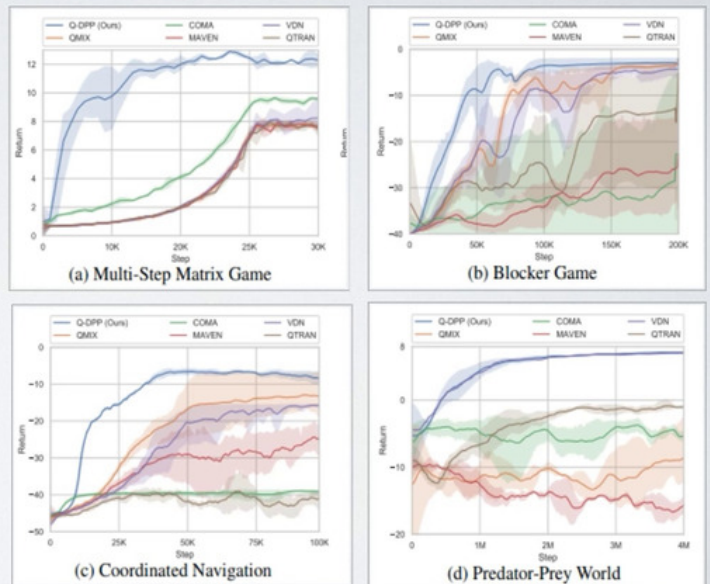


Figure 3: Multi-agent cooperative tasks. The size of the ground set for each task is a) 176, b) 420, c) 720, d) 3920.

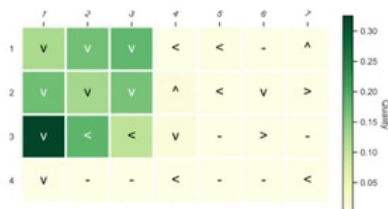
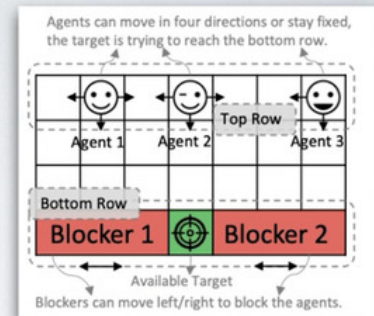


Some Insights: What is the Learned Diverse Behavioural Model

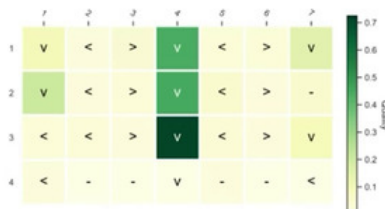
• The acquired diverse behaviours on Blocker Game.

- The optimal action of one agent does not depend on the others
- The behaviours acquired by different agents are orthogonal.
- Q-DPP meets the decentralisability assumption naturally.

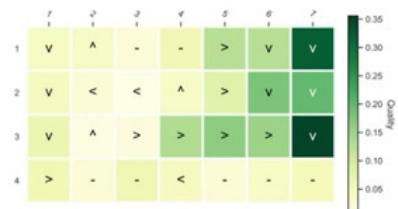
$$\arg \max_a Q^\pi(s, a) = \begin{bmatrix} \arg \max_{a_1} Q_1(s_1, a_1) \\ \vdots \\ \arg \max_{a_N} Q_N(s_N, a_N) \end{bmatrix}$$



(a) Agent 1.



(b) Agent 2.



(c) Agent 3.

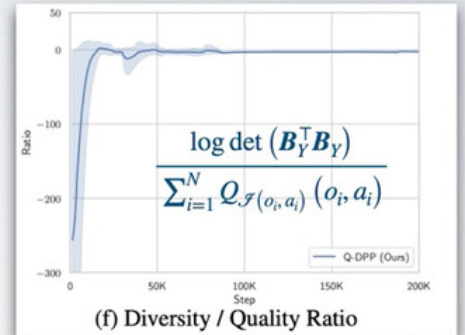
Some Insights: Quality vs Diversity

- The importance of diversity decrease with training.
- Q-DPP generalizes VDN, QMIX, QTRAN.

$$Q^\pi(o, \mathbf{a}^*) = \sum_{i=1}^N Q_{\mathcal{F}(o_i, \mathbf{a}_i^*)}(o_i, \mathbf{a}_i^*) + \log \det(\mathbf{B}_Y^\top \mathbf{B}_Y = \mathbf{I})$$

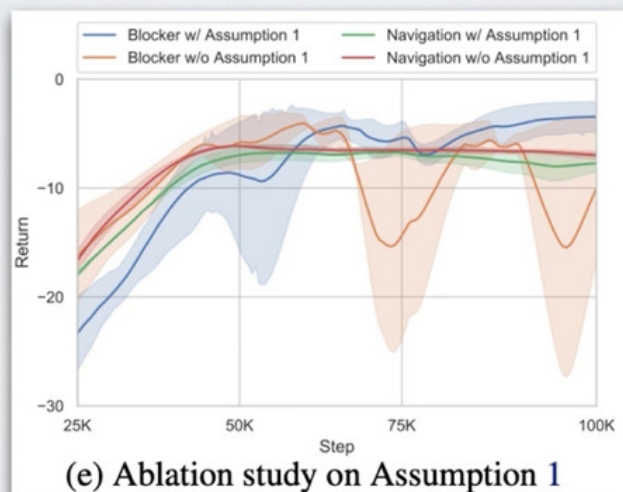
- Recovers VDN: $Q^\pi(o, \mathbf{a}^*) = \sum_{i=1}^N Q_{\mathcal{F}(o_i, \mathbf{a}_i^*)}(o_i, \mathbf{a}_i^*)$
- Recovers QMIX: $\frac{\partial Q^\pi(o, \mathbf{a}^*)}{\partial Q_{\mathcal{F}(o_i, \mathbf{a}_i^*)}(o_i, \mathbf{a}_i^*)} = 1 \geq 0, \quad \forall \mathcal{F}(o_i, \mathbf{a}_i^*) \in \mathcal{N}$
- Recovers QTRAN: $\sum_{i=1}^N Q_i(o_i, a_i) - Q^\pi(o, \mathbf{a}) + V(o) = \begin{cases} 0 & \mathbf{a} = \mathbf{a}^* \\ \geq 0 & \mathbf{a} \neq \mathbf{a}^* \end{cases}$

- Diversity is defined from the perspective of maximising the reward. Agents are still allowed to take the same actions in a task.



Some Insights: the Effect of Assumption 1

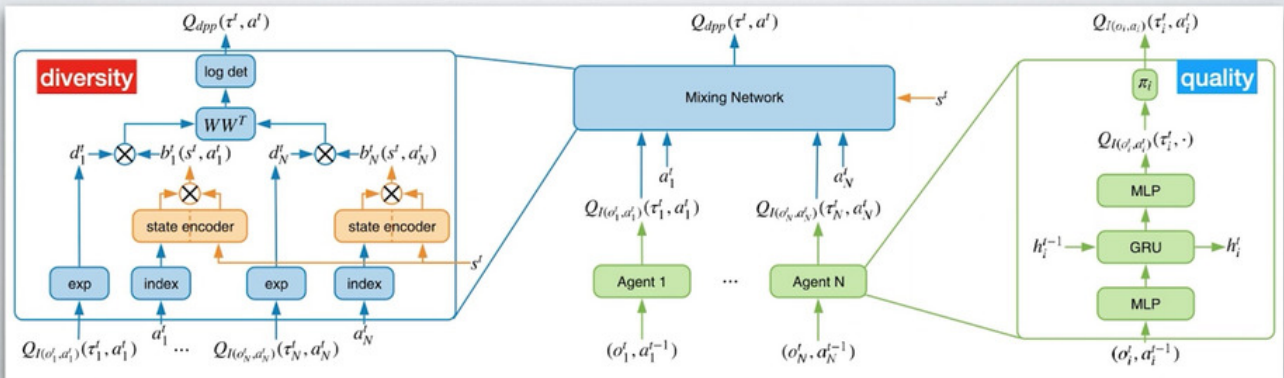
- Assumption 1 reduces the variance of the rewards. It is because it mitigates the over-estimation of the samples from the biased Q-DPP sampler, which helps agents avoid getting over-confident, as a result, trapped in the bad local optima.



Q-DPP on Continuous State Space

• Proof of concept on the continuous-state tasks.

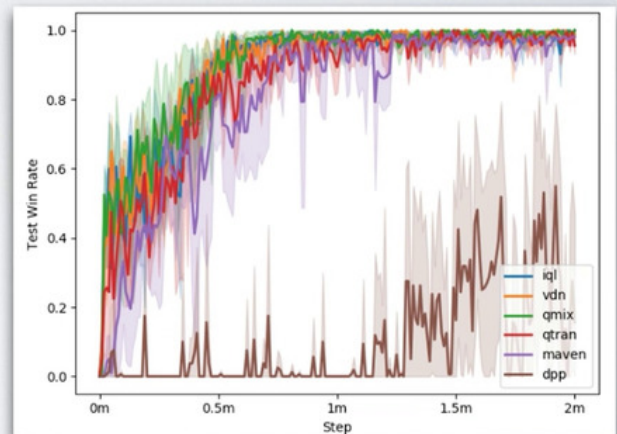
- Both the quality term and the diversity term can be modelled by a neural network.
- The quality term that models $Q_i(s, a_i)$ is to the most right.
- The diversity term is modelled by a directional network and a norm network separately.

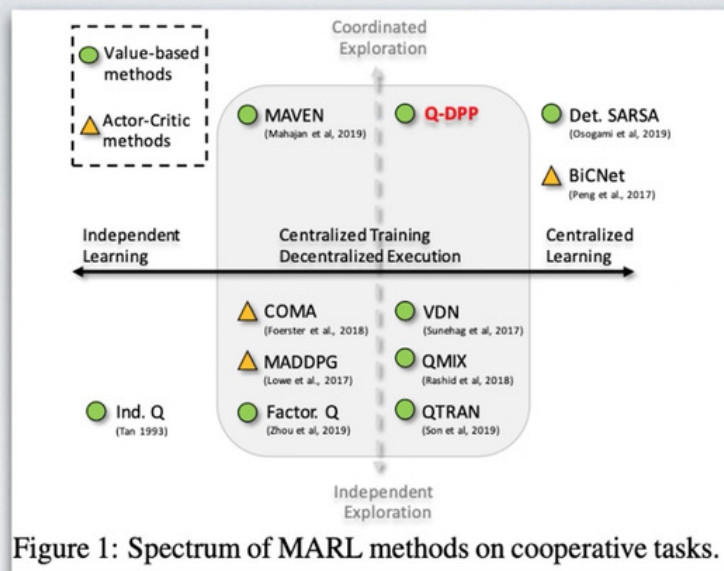


Q-DPP on StarCraft II Micro-management

- **Proof of concept on the continuous-state tasks.**

- How to model continuous set via DPPs are still under development.
- We believe a full treatment of the idea of **Depp Q-DPP** needs substantial future work.





Thank you!