

# CS4400

# DEEP REINFORCEMENT LEARNING

## Lecture 11: Advanced MARL

Wendelin Böhmer

`<j.w.bohmer@tudelft.nl>`



11th of January 2024

# Content of this lecture



- 11.1 Value factorization
- 11.2 Relative overgeneralization
- 11.3 Communication

11.1

# Advanced MARL

## Value factorization

- Cooperative tasks allow centralized objectives
  - with decentralized decision policies

- Decentralizable value factorization  $Q_{\theta}(\tau_t, \mathbf{a}) := \sum_{i=1}^N q_{\theta}^i(\tau_t^i, a^i)$ 
  - called *value decomposition networks* (VDN)
  - *utilities*  $q_{\theta}^i$  do not fulfill Bellman equation

$$\mathbf{a}_t^* = \arg \max_{\mathbf{a} \in \mathcal{A}} Q_{\theta}(\tau_t, \mathbf{a}) = \left\{ \arg \max_{a^i \in \mathcal{A}^i} q_{\theta}^i(\tau_t^i, a^i) \right\}_{i=1}^N$$

- Optimized with centralized (double) DQN loss

$$\mathcal{L}_{[\theta]}^{\text{VDN}} := \mathbb{E}_{\mathcal{D}} \left[ \left( r_t + \gamma Q_{\theta'}(\tau_{t+1}, \arg \max_{\mathbf{a} \in \mathcal{A}} Q_{\theta}(\tau_{t+1}, \mathbf{a})) - Q_{\theta}(\tau_t, \mathbf{a}_t) \right)^2 \right]$$



- VDN factorization restricts function class of  $Q_\theta$
- Any monotonic value mixture is decentralizable:  $\frac{\partial Q_\theta(\tau_t, \mathbf{a})}{\partial q_\theta^i(\tau_t^i, a^i)} \geq 0$ 
  - e.g. VDN is monotonic  $\frac{\partial \sum_{j=1}^N q_\theta^j(\tau_t^j, a^j)}{\partial q_\theta^i(\tau_t^i, a^i)} = 1$



- VDN factorization restricts function class of  $Q_\theta$
- Any monotonic value mixture is decentralizable:  $\frac{\partial Q_\theta(\tau_t, \mathbf{a})}{\partial q_\theta^i(\tau_t^i, a^i)} \geq 0$ 
  - e.g. VDN is monotonic  $\frac{\partial \sum_{j=1}^N q_\theta^j(\tau_t^j, a^j)}{\partial q_\theta^i(\tau_t^i, a^i)} = 1$
- QMIX: learn centralized monotonic mixture network  $f_\phi$ 
  - $f_\phi$  not required during execution
  - $f_\phi$  yields different mixing for different states  $s_t$

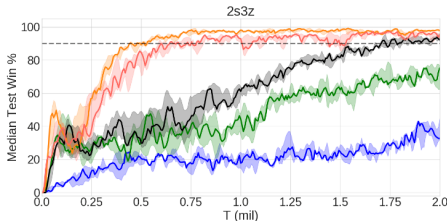
$$Q_{\theta\phi}(s_t, \tau_t, \mathbf{a}) := f_\phi(s_t, q_\theta^1(\tau_t^1, a^1), \dots, q_\theta^N(\tau_t^N, a^N)) =: f_\phi(s_t, \mathbf{q}_\theta(\tau_t, \mathbf{a}))$$

- *Hyper-network*  $f_\phi$  with non-negative weights, e.g.:

$$f_\phi(\mathbf{s}, \mathbf{q}) = \mathbf{w}_{(\mathbf{s})}^{1\top} \sigma(\mathbf{W}_{(\mathbf{s})}^2 \mathbf{q} + \mathbf{w}_{(\mathbf{s})}^3) + w_{(\mathbf{s})}^4, \quad \mathbf{w}_{(\mathbf{s})}^k = \underbrace{|\mathbf{A}_2^k|}_{\text{hyper-net parameters}} \underbrace{\bar{\sigma}(\mathbf{A}_1^k \mathbf{s} + \mathbf{b}_1^k)}_{\phi} + \underbrace{|\mathbf{b}_2^k|}_{\phi}$$

QMIX by [Rashid et al. \(2018, 2020b\)](#)

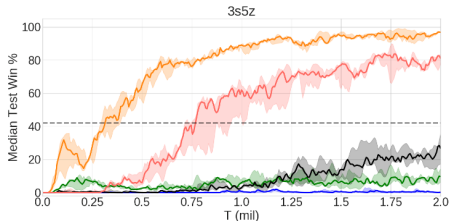
- Play against build-in AI
- Partial observability
- Varying units available
  - here: stalkers and zealots



QMIX

VDN

QTRAN



IQL

COMA

StarCraft II multi-agent challenge (SMAC, [Samvelyan et al., 2019](#)); QTRAN ([Son et al., 2019](#)); results from ([Rashid et al., 2020b](#))



- Centralized value functions can be *factorized*
- Factorization must be monotonic in per-agent utilities
- Maximizing utilities independently maximizes joint value
- QMIX extends VDN by using non-negative hyper-networks
- QTRAN defines decentrability as constraints, but is instable

### Learning Objectives

LO11.1: Explain monotonic value factorization in VDN and QMIX



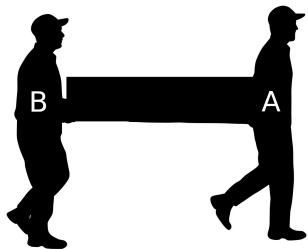
11.2

# **Advanced MARL**

## Relative overgeneralization



- Why is independent learning or value factorization problematic?
  - e.g. single-state IQL  $q^i(a^i; \pi^{-i}) := \mathbb{E}[Q^\pi(a) \mid a^{-i} \sim \pi^{-i}(\cdot)]$
- What are the IQL values if all agents explore randomly?
  - team needs to transport a fragile box (reward +1)
  - the box falls (punishment -1) if A is too far from B
  - will a team of IQL agents learn an optimal policy?



		A				
		←	∅	→		
B	←	0	-1	-1		
	∅	0	0	-1		
	→	0	0	+1		
$q_{\text{explore}}^A$		?				

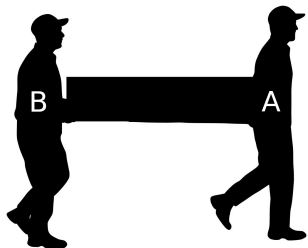
image source: [pxhere.com](https://pxhere.com)



# core concept: Relative overgeneralization



- Why is independent learning or value factorization problematic?
  - e.g. single-state IQL  $q^i(a^i; \pi^{-i}) := \mathbb{E}[Q^\pi(a) \mid a^{-i} \sim \pi^{-i}(\cdot)]$
- What are the IQL values if all agents explore randomly?
  - agent A learns **incorrectly** to go *backwards*



		A			$q_{\text{exp}}^B$	$q_{\text{gre}}^B$
		$\leftarrow$	$\emptyset$	$\rightarrow$		
B	$\leftarrow$	0	-1	-1	?	
	$\emptyset$	0	0	-1		
	$\rightarrow$	0	0	+1		
$q_{\text{explore}}^A$		0	$-\frac{1}{3}$	$-\frac{1}{3}$		

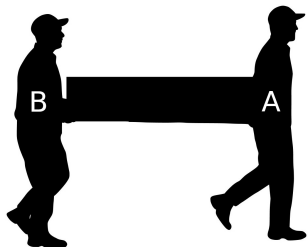
this exploration effect on independent values is called *relative overgeneralization* (Panait et al., 2006); image source: [pxhere.com](http://pxhere.com)



# core concept: Relative overgeneralization



- Why is independent learning or value factorization problematic?
  - e.g. single-state IQL  $q^i(a^i; \pi^{-i}) := \mathbb{E}[Q^\pi(a) \mid a^{-i} \sim \pi^{-i}(\cdot)]$
- What are the IQL values if all agents explore randomly?
  - agent A learns **incorrectly** to go *backwards*
  - agent B learns **correctly** to go *forwards*
- Will the team learn the task when they act more greedy?



		A			$q_{\text{exp}}^B$	$q_{\text{gre}}^B$
		$\leftarrow$	$\emptyset$	$\rightarrow$		
B	$\leftarrow$	0	-1	-1	$-\frac{2}{3}$	?
	$\emptyset$	0	0	-1	$-\frac{1}{3}$	
	$\rightarrow$	0	0	+1	$\frac{1}{3}$	
$q_{\text{explore}}^A$		0	$-\frac{1}{3}$	$-\frac{1}{3}$		
$q_{\text{greedy}}^A$			?			

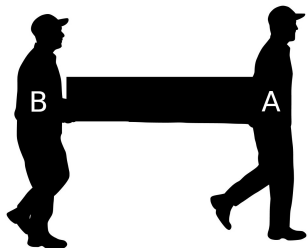
this exploration effect on independent values is called *relative overgeneralization* (Panait et al., 2006); image source: [pxhere.com](https://pxhere.com)



# core concept: Relative overgeneralization



- Why is independent learning or value factorization problematic?
  - e.g. single-state IQL  $q^i(a^i; \pi^{-i}) := \mathbb{E}[Q^\pi(a) \mid a^{-i} \sim \pi^{-i}(\cdot)]$
- What are the IQL values if all agents explore randomly?
  - agent A learns **incorrectly** to go *backwards*
  - agent B learns **correctly** to go *forwards*
- Will the team learn the task when they act more greedy?
  - $q^A(a_{\rightarrow}^A) = \pi^B(a_{\rightarrow}^B) - (1 - \pi^B(a_{\rightarrow}^B))$ ,  $\pi^B(a_{\rightarrow}^B) > \frac{1}{2} \Rightarrow q^A(a_{\rightarrow}^B) > 0$



		A			$q_{\text{exp}}^B$	$q_{\text{gre}}^B$
		$\leftarrow$	$\emptyset$	$\rightarrow$		
B	$\leftarrow$	0	-1	-1	$-\frac{2}{3}$	0
	$\emptyset$	0	0	-1	$-\frac{1}{3}$	0
	$\rightarrow$	0	0	+1	$\frac{1}{3}$	+1
$q_{\text{explore}}^A$		0	$-\frac{1}{3}$	$-\frac{1}{3}$		
$q_{\text{greedy}}^A$		0	0	+1		

this exploration effect on independent values is called *relative overgeneralization* (Panait et al., 2006); image source: [pxhere.com](http://pxhere.com)



## Question: unlearnable games



- Create a reward matrix for a one-state 2 player cooperative game, which **cannot** be learned by IQL agents





## Question: unlearnable games



- Create a reward matrix for a one-state 2 player cooperative game, which **cannot** be learned by IQL agents
- for example the following table
  - note that an unfactored value  $Q^\pi(a)$  could solve the task

		A			$q_{\text{exp}}^B$
		$a_1^A$	$a_2^A$	$a_3^A$	
B	$a_1^B$	0	0	-2	$-\frac{2}{3}$
	$a_2^B$	0	0	-2	$-\frac{2}{3}$
	$a_3^B$	-2	-2	+1	-1
$q_{\text{explore}}^A$		$-\frac{2}{3}$	$-\frac{2}{3}$	-1	



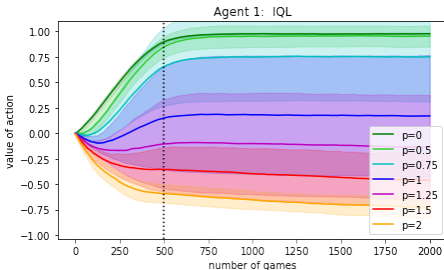


## Question: unlearnable games



- Create a reward matrix for a one-state 2 player cooperative game, which **cannot** be learned by IQL agents
- for example the following table
  - note that an unfactored value  $Q^\pi(a)$  could solve the task
  - punishment  $p$  only increases *probability* of failure

		A			$q_{\text{exp}}^B$
		$a_1^A$	$a_2^A$	$a_3^A$	
B	$a_1^B$	0	0	$-p$	$-\frac{p}{3}$
	$a_2^B$	0	0	$-p$	$-\frac{p}{3}$
	$a_3^B$	$-p$	$-p$	$+1$	$\frac{1-2p}{3}$
$q_{\text{explore}}^A$		$-\frac{p}{3}$	$-\frac{p}{3}$	$\frac{1-2p}{3}$	



assignment sheet 4



## 11.2 Optimistic return methods

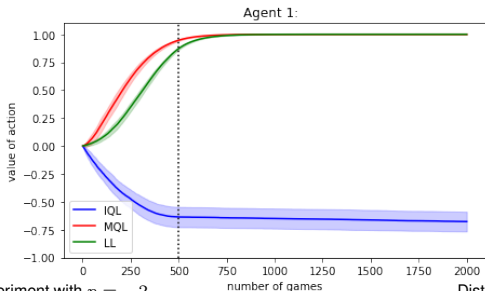


- IQL always learn when partners behave optimal

$$\text{TD}(\tau_{t+1}^i, r_t^i) = r_t^i + \gamma \max_{a'^i} q_{\theta^i}^i(\tau_{t+1}^i, a'^i) - q_{\theta^i}^i(\tau_t^i, a_t^i)$$

- Distributed Q-learning pretends others play best

- use **best past reward**:  $\text{TD}(\tau_{t+1}^i, \max\{r_{t'}^i \mid s_{t'}=s_t, a_{t'}=a_t\})$
- or **ignore negative TD errors**:  $\max(0, \text{TD}(\tau_{t+1}^i, r_t^i))$
- both fail in stochastic environments



$\epsilon$ -greedy experiment with  $p = -2$

Distributed Q-learning (Lauer and Riedmiller, 2000)



- **IQL** always learn when partners behave optimal

$$\text{TD}(\tau_{t+1}^i, r_t^i) = r_t^i + \gamma \max_{a'^i} q_{\theta^i}^i(\tau_{t+1}^i, a'^i) - q_{\theta^i}^i(\tau_t^i, a_t^i)$$

- Distributed Q-learning pretends others play best

- use **best past reward**:  $\text{TD}(\tau_{t+1}^i, \max\{r_{t'}^i \mid s_{t'}=s_t, a_{t'}=a_t\})$
- or **ignore negative TD errors**:  $\max(0, \text{TD}(\tau_{t+1}^i, r_t^i))$
- both fail in stochastic environments

- Lenient learning reduces slowly to IQL

- **ignore negative TD errors** only *sometimes*
- decreases ignoring probability  $\eta(s, a)$  with visitations
- deep version must remember  $\eta(s, a)$  in hash table

Lenient learning by [Panait et al. \(2008\)](#), deep lenient learning by [Wei and Luke \(2016\)](#)

- **IQL** always learn when partners behave optimal

$$\text{TD}(\tau_{t+1}^i, r_t^i) = r_t^i + \gamma \max_{a'^i} q_{\theta^i}^i(\tau_{t+1}^i, a'^i) - q_{\theta^i}^i(\tau_t^i, a_t^i)$$

- Distributed Q-learning pretends others play best

- use **best past reward**:  $\text{TD}(\tau_{t+1}^i, \max\{r_{t'}^i \mid s_{t'}=s_t, a_{t'}=a_t\})$
- or **ignore negative TD errors**:  $\max(0, \text{TD}(\tau_{t+1}^i, r_t^i))$
- both fail in stochastic environments

- Lenient learning reduces slowly to IQL

- **ignore negative TD errors** only *sometimes*
- decreases ignoring probability  $\eta(s, a)$  with visitations
- deep version must remember  $\eta(s, a)$  in hash table

- Weighted-QMIX weights some outcomes with  $\alpha \in [0, 1]$

- QMIX with additional centralized Q-value function
- weights TD-errors with  $\alpha$  if greedy policies disagree

wighted QMIX by [Rashid et al. \(2020b\)](#), other approaches sample *similar tasks* optimistically (e.g. [Gupta et al., 2021](#))

- Unfactored value functions have huge action space
- Use higher order functions can coordinate more agents

$$Q^\pi(\tau_t, \mathbf{a}) := \underbrace{\sum_{i \in \mathcal{E}^1} q^i(\tau_t^i, a^i) + \sum_{(i,j) \in \mathcal{E}^2} q^{ij}(\tau_t^{ij}, \mathbf{a}^{ij}) + \sum_{(i,j,k) \in \mathcal{E}^3} \dots}_{\text{COORDINATION GRAPH}} \quad \text{VDN}$$

max-plus algorithm see [Pearl \(1988\)](#), CG introduced by [Guestrin et al. \(2002\)](#), tabular CG used in [Kok and Vlassis \(2006\)](#)

- Unfactored value functions have huge action space
- Use higher order functions can coordinate more agents

$$Q^\pi(\tau_t, \mathbf{a}) := \underbrace{\sum_{i \in \mathcal{E}^1} \overbrace{q^i(\tau_t^i, a^i)}^{\text{VDN}} + \sum_{(i,j) \in \mathcal{E}^2} q^{ij}(\tau_t^{ij}, \mathbf{a}^{ij}) + \sum_{(i,j,k) \in \mathcal{E}^3} \dots}_{\text{COORDINATION GRAPH}}$$

- Maximum no longer decentralizable
  - e.g. *coordination graph* (CG)  $\mathcal{G} := \langle \mathcal{E}^1, \mathcal{E}^2 \rangle$
  - all agents are nodes:  $\mathcal{E}^1 := \{i\}_{i=1}^N$
  - pairwise coordination:  $\mathcal{E}^2 \subseteq \{(i, j) | 1 \leq i \leq N, 1 \leq j \leq N\}$
  - maximum can be computed with *max-plus algorithm*
  - computed using multiple message passes between agents

max-plus algorithm see [Pearl \(1988\)](#), CG introduced by [Guestrin et al. \(2002\)](#), tabular CG used in [Kok and Vlassis \(2006\)](#)

- Deep coordination graphs (DCG) implements pairwise edges

$$Q_{\theta\phi\psi\varphi}^{\text{DCG}}(s_t, \tau_t, \mathbf{a}) := q_{\varphi}^0(s_t) + \sum_{i \in \mathcal{E}^1} q^i(\tau_t^i, a^i) + \sum_{(i,j) \in \mathcal{E}^2} q^{ij}(\tau_t^{ij}, \mathbf{a}^{ij})$$

- Required extensive engineering to work with neural nets

- shared history encoding:  $\mathbf{h}_t^i := h_{\psi}(\tau_t^i) = h_{\psi}(\mathbf{h}_{t-1}^i, o_t^i, a_{t-1}^i)$
- shared utility functions in  $\mathcal{E}^1$ :  $q^i(\tau_t^i, a^i) := q_{\theta}^v(\mathbf{h}_t^i, a^i)$
- symmetrized shared payoff functions in  $\mathcal{E}^2$ :

$$q^{ij}(\tau_t^{ij}, \mathbf{a}^{ij}) := \frac{1}{2} (q_{\phi}^e(\mathbf{h}_t^i, \mathbf{h}_t^j, a^i, a^j) + q_{\phi}^e(\mathbf{h}_t^j, \mathbf{h}_t^i, a^j, a^i))$$

- $K$ -rank approximation of payoff matrices:

$$q_{\phi}^e(\mathbf{h}_t^i, \mathbf{h}_t^j, a^i, a^j) := \sum_{k=1}^K g_{\phi}^k(\mathbf{h}_t^i, \mathbf{h}_t^j, a^i) \bar{g}_{\phi}^k(\mathbf{h}_t^i, \mathbf{h}_t^j, a^j)$$

- state-dependent bias  $q_{\varphi}^0(s_t)$

[Böhmer et al. \(2020\)](#) introduce DCG, and [Castellini et al. \(2019\)](#) explore higher order factors

## 11.2 Comparison: predator-prey



- 10x10 gridworld
- 8 agents, 8 prey
- reward +10/prey
- punishment -2/attempt
- catch removes agents+prey

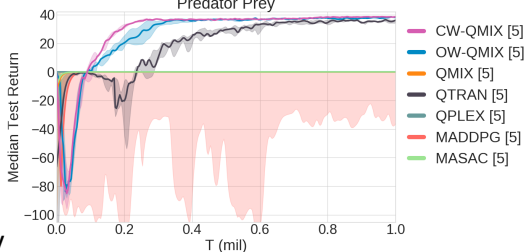
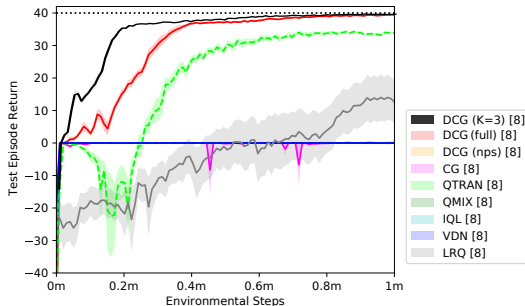


image and upper results from [Böhmer et al. \(2020\)](#), lower results from ([Rashid et al., 2020a](#))



- Random exploration can cause RO in independent learners
- When expected return is low only because others explore
- Can be counteracted by down-weighting errors with low returns
- Higher-order value factorization can represent joint-Q-values
- DCG requires message passing and many stabilization tricks

## Learning Objectives

LO11.2: Derive when a game exhibits relative overgeneralization

LO11.3: Explain what can be done against relative overgeneralization



11.3

## **Advanced MARL** Communication



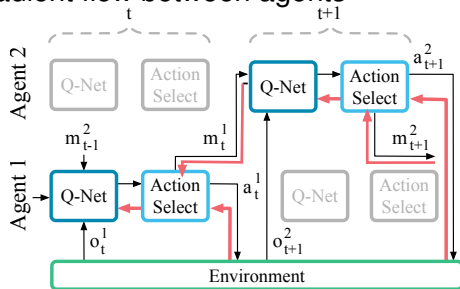
- Coordination graphs use communication for message passing
  - communication protocol fixed by max-plus algorithm
  - can we learn a useful protocol between agents?
- Communication almost always possible
  - *implicitly* by manipulating the environment
  - *explicitly* by using “cheap-talk” channels
- Not every game benefits from communication
  - no incentive to help opponent in zero-sum games
  - messages in general-sum games can be antagonistic
  - fully observable cooperation games do not need communication
  - e.g. communicate all observations for centralized agent



# core concept: Differentiable communication



- Learn communication implicitly as part of architecture
  - requires additional *cheap-talk* channels
  - messages  $m_t^i$  at time  $t$  are inputs of time  $t + 1$
- Continuous messages allow gradient flow between agents



e.g. Wang et al. (2018) use graph neural networks for message passing, figure modified from Förster et al. (2016)



# core concept: Differentiable communication



- Learn communication implicitly as part of architecture
  - requires additional *cheap-talk* channels
  - messages  $m_t^i$  at time  $t$  are inputs of time  $t + 1$
- Continuous messages allow gradient flow between agents
- No gradient flow through sampled messages!
  - $m_t^i \in \{1, \dots, M\}$
  - $m_t^i \sim \mu_\theta^i(\cdot | \tau_t^i)$

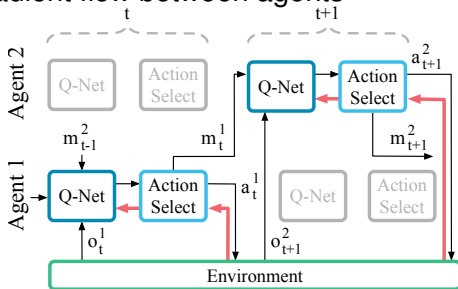


figure from Förster et al. (2016), sending (non-differentiable) messages between RNN-agents



- Learn communication implicitly as part of architecture
  - requires additional *cheap-talk* channels
  - messages  $m_t^i$  at time  $t$  are inputs of time  $t + 1$
- Continuous messages allow gradient flow between agents
- No gradient flow through sampled messages!
  - $m_t^i \in \{1, \dots, M\}$ ,  $\epsilon \sim \mathcal{N}(\cdot|0, 1)$
  - $m_t^i = f_\theta(\tau_t^i, \epsilon) \sim \mu_\theta^i(\cdot|\tau_t^i)$
- Use reparametrization trick
  - gradient flows through  $f_\theta$
  - $\nabla_\theta \mathbb{E} \left[ g(m_t^i) \mid m_t^i \sim \mu_\theta^i(\cdot|\tau_t^i) \right] = \mathbb{E} \left[ \nabla_\theta g(f_\theta(\tau_t^i, \epsilon)) \mid \epsilon \sim \mathcal{N}(\cdot|0, 1) \right]$

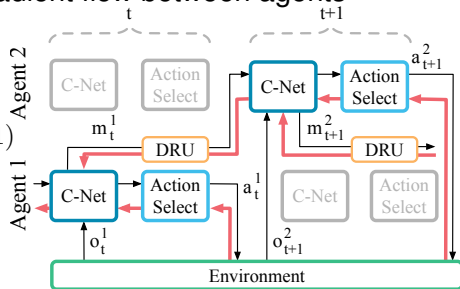
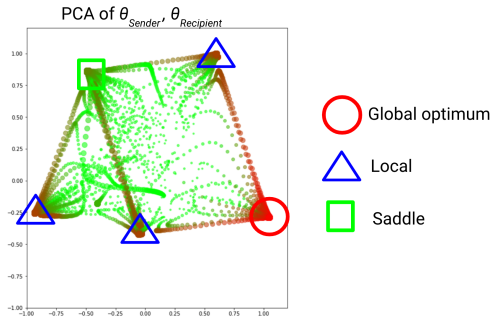


figure with reparametrization modules (DRU) from [Förster et al. \(2016\)](#), see also [Jang et al. \(2017\)](#) for differentiable sampling



- Can we learn to communicate in cooperative games?
  - provide some communication actions, e.g.  $m_t^i \in \{0, 1\}$
  - negotiate meaning of  $m_t^i$  during training
- Established conventions are hard to change, e.g.
  - A sees  $\triangle \Rightarrow m^A = 1$
  - $m^A = 1 \Rightarrow$  B does  $\triangle$
- Many Nash-equilibria!

		B does		
		○	△	▽
A sees	○	+1	-1	-1
	△	-1	$+\frac{1}{2}$	-1
	▽	-1	-1	$+\frac{1}{2}$



example and PCA plot by Jakob Förster



- Communication can help to coordinate agents
- Not all games benefit from communication
- Differentiable communication in centralized training
- Discrete channels must use the reparameterization trick
- Learned communication protocols have many Nash equilibria

### Learning Objectives

LO11.4: Explain in which games communication is beneficial

LO11.5: Explain differentiable communication and how it can fail

- Next lecture: **applied RL**!
- Questions? Ask them here: [answers.ewi.tudelft.nl](https://answers.ewi.tudelft.nl)

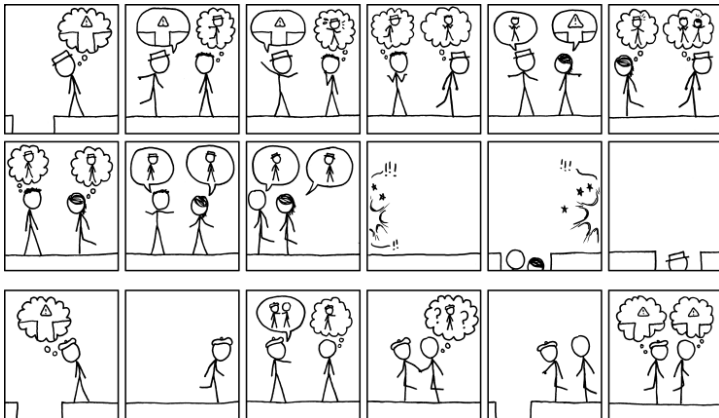


image source: [xkcd.com](https://xkcd.com)



- Wendelin Böhmer, Vitaly Kurin, and Shimon Whiteson. Deep coordination graphs. In *Proceedings of Machine Learning and Systems (ICML)*, pages 2611–2622, 2020. URL <https://arxiv.org/abs/1910.00091>.
- Jacopo Castellini, Frans A. Oliehoek, Rahul Savani, and Shimon Whiteson. The representational capacity of action-value networks for multi-agent reinforcement learning. In *Proceedings of the 18th International Conference on Autonomous Agents and MultiAgent Systems, AAMAS '19*, pages 1862–1864, 2019. URL <http://www.ifaamas.org/Proceedings/aamas2019/pdfs/pl1862.pdf>.
- Jakob Förster, Ioannis Alexandros Assael, Nando de Freitas, and Shimon Whiteson. Learning to communicate with deep multi-agent reinforcement learning. In *Advances in Neural Information Processing Systems 29*, pages 2137–2145. 2016. URL <http://papers.nips.cc/paper/6042-learning-to-communicate-with-deep-multi-agent-reinforcement-learning.pdf>.
- Carlos Guestrin, Michail Lagoudakis, and Ronald Parr. Coordinated reinforcement learning. In *ICML*, volume 2, pages 227–234, 2002.
- Tarun Gupta, Anuj Mahajan, Bei Peng, Wendelin Böhmer, and Shimon Whiteson. UneVEN: Universal value exploration for multi-agent reinforcement learning. In *Proceedings of the International Conference on Machine Learning (ICML)*. 2021. URL <https://arxiv.org/abs/2010.02974>.
- Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparametrization with gumbel-softmax. In *Proceedings International Conference on Learning Representations 2017*, 2017. URL <https://openreview.net/pdf?id=rkE3y85ee>.
- Jelle R Kok and Nikos Vlassis. Collaborative multiagent reinforcement learning by payoff propagation. *Journal of Machine Learning Research*, 7(Sep):1789–1828, 2006.
- Martin Lauer and Martin Riedmiller. An algorithm for distributed reinforcement learning in cooperative multi-agent systems. In *In Proceedings of the Seventeenth International Conference on Machine Learning*, pages 535–542. Morgan Kaufmann, 2000.
- Liviu Panait, Sean Luke, and R. Paul Wiegand. Biasing coevolutionary search for optimal multiagent behaviors. *IEEE Transactions on Evolutionary Computation*, 10(6):629–645, 2006.



- Liviu Panait, Karl Tuyls, and Sean Luke. Theoretical advantages of lenient learners: An evolutionary game theoretic perspective. *The Journal of Machine Learning Research*, 9:423–457, 2008.
- Judea Pearl. *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1988. ISBN 0-934613-73-7.
- Tabish Rashid, Mikayel Samvelyan, Christian Schröder de Witt, Gregory Farquhar, Jakob N. Förster, and Shimon Whiteson. QMIX: monotonic value function factorisation for deep multi-agent reinforcement learning. In *International Conference on Machine Learning (ICML)*, pages 4292–4301, 2018.
- Tabish Rashid, Gregory Farquhar, Bei Peng, and Shimon Whiteson. Weighted qmix: Expanding monotonic value function factorisation. In *Advances in Neural Information Processing Systems*, 2020a. URL <https://arxiv.org/abs/2006.10800>.
- Tabish Rashid, Mikayel Samvelyan, Christian Schröder de Witt, Gregory Farquhar, Jakob N. Förster, and Shimon Whiteson. Monotonic value function factorisation for deep multi-agent reinforcement learning. *Journal of Machine Learning Research (JMLR)*, 2020b. URL <https://arxiv.org/abs/2003.08839>.
- Mikayel Samvelyan, Tabish Rashid, Christian Schroeder de Witt, Gregory Farquhar, Nantas Nardelli, Tim G. J. Rudner, Chia-Man Hung, Philip H. S. Torr, Jakob Förster, and Shimon Whiteson. The StarCraft Multi-Agent Challenge. *CoRR*, abs/1902.04043, 2019.
- Kyunghwan Son, Daewoo Kim, Wan Ju Kang, David Earl Hostallero, and Yung Yi. QTRAN: Learning to factorize with transformation for cooperative multi-agent reinforcement learning. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 5887–5896, 2019. URL <http://proceedings.mlr.press/v97/son19a.html>.
- Peter Sunehag, Guy Lever, Audrunas Gruslys, Wojciech Marian Czarnecki, Vinicius Zambaldi, Max Jaderberg, Marc Lanctot, Nicolas Sonnerat, Joel Z. Leibo, Karl Tuyls, and Thore Graepel. Value-decomposition networks for cooperative multi-agent learning based on team reward. In *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems (AAMAS)*, pages 2085–2087, 2018.



Tingwu Wang, Renjie Liao, Jimmy Ba, and Sanja Fidler. Nervenet: Learning structured policy with graph neural networks. In *ICLR*, 2018. URL <https://openreview.net/forum?id=S1sqHMZCb>.

Ermo Wei and Sean Luke. Lenient learning in independent-learner stochastic cooperative games. *The Journal of Machine Learning Research*, 17(1):2914–2955, 2016.