

# QPLEX: Duplex Dueling Multi-Agent Q-Learning

Jianhao Wang<sup>\*</sup>, Zhizhou Ren<sup>\*</sup>, Terry Liu, Yang Yu, Chongjie Zhang

Tsinghua University; Polixir Technologies

<sup>\*</sup>equal contribution



Machine Intelligence Group



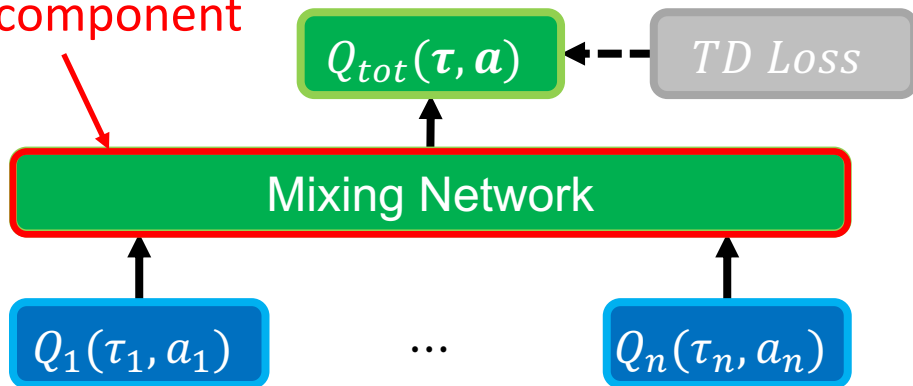
清华大学  
Tsinghua University

交叉信息研究院  
Institute for Interdisciplinary Information Sciences

# Value Function Factorization Methods

- Paradigm: centralized training with decentralized execution

Core component



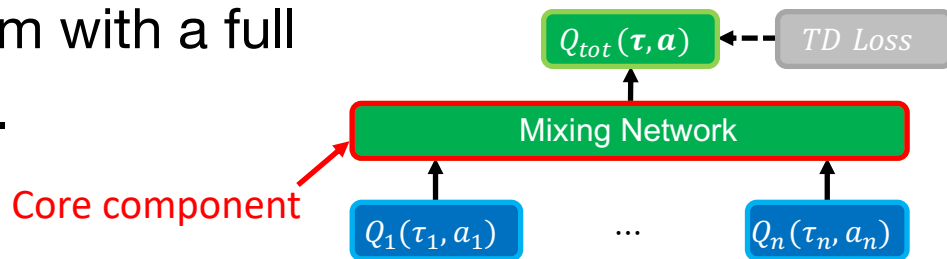
$$\mathcal{L}(\theta) = \mathbb{E} \left[ \left( r + \gamma V(\tau'; \theta^-) - Q(\tau, a; \theta) \right)^2 \right]$$

$$V(\tau'; \theta^-) = \max_{a'} Q(\tau', a'; \theta^-)$$



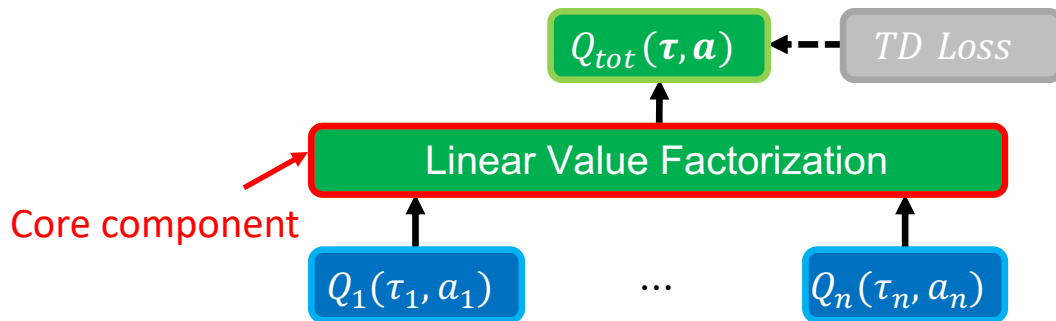
# Individual-Global-Maximization (IGM) Principle

- Paradigm: centralized training with decentralized execution
- Individual-Global Maximization (IGM) Principle
  - Consistent greedy action selection between joint and individuals
  - $\operatorname{argmax}_a Q_{tot}(\boldsymbol{\tau}, \mathbf{a}) = (\operatorname{argmax}_{a_1} Q_1(\tau_1, a_1), \dots, \operatorname{argmax}_{a_n} Q_n(\tau_n, a_n))$
  - **Open problem:** How to design an effective and scalable multi-agent Q-learning algorithm with a full realization of the IGM principle.



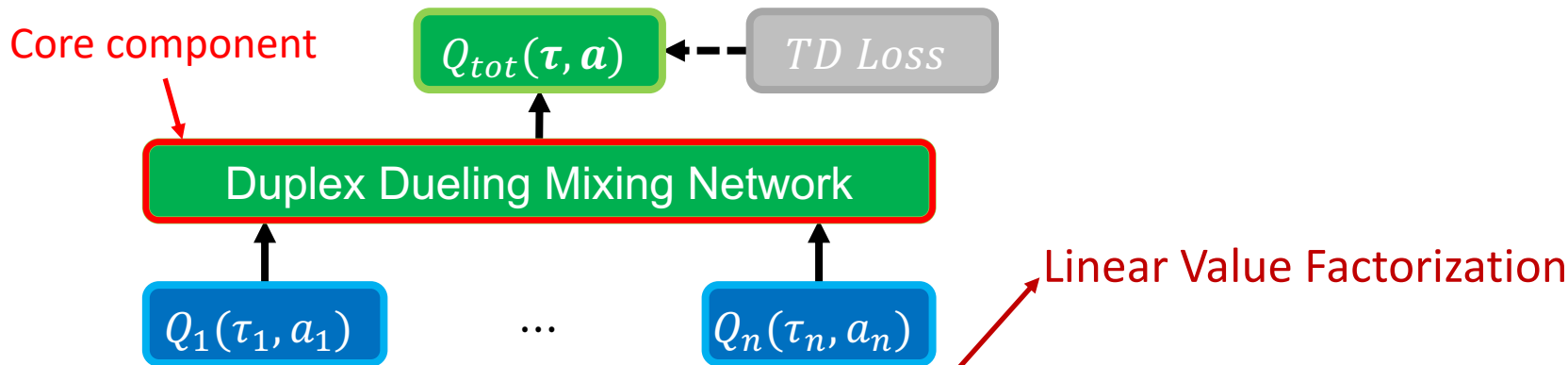
# Previous Work: Value Decomposition Networks (VDN)

- VDN:  $Q_{tot}(\boldsymbol{\tau}, \boldsymbol{a}) = \sum_i Q_i(\tau_i, a_i)$
- Sufficient for IGM constraint
  - $\operatorname{argmax}_{\boldsymbol{a}} Q_{tot}(\boldsymbol{\tau}, \boldsymbol{a}) = (\operatorname{argmax}_{a_1} Q_1(\tau_1, a_1), \dots, \operatorname{argmax}_{a_n} Q_n(\tau_n, a_n))$
- Excellent scalability
- Cons:
  - Not necessary for IGM
  - Limited expressiveness



# QPLEX: Duplex Dueling Multi-Agent Q-Learning

- Paradigm: centralized training with decentralized execution

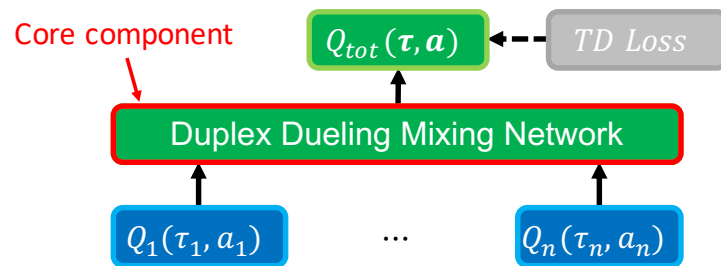


- QPLEX :  $Q_{tot}(\tau, a) = \sum_i Q_i(\tau, a_i) + \sum_{i=1}^n (\lambda_i(\tau, a) - 1) A_i(\tau, a_i)$ 
  - Strong representation capacity
  - Easily realized and learned by neural networks

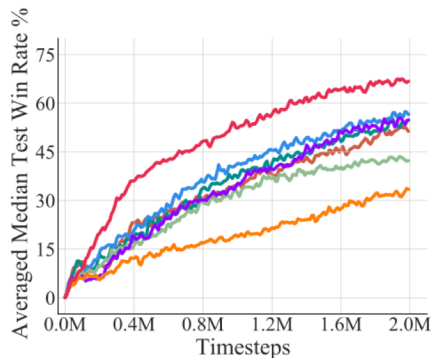


# QPLEX: Duplex Dueling Multi-Agent Q-Learning

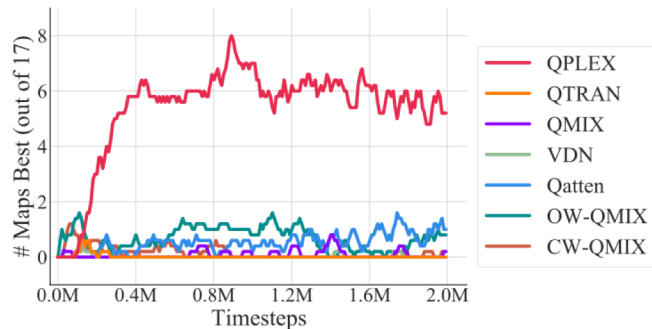
Theorem: *The joint action-value function class that QPLEX can realize is **equivalent** to what is induced by the IGM principle.*



# StarCraft II Benchmark: Online Learning

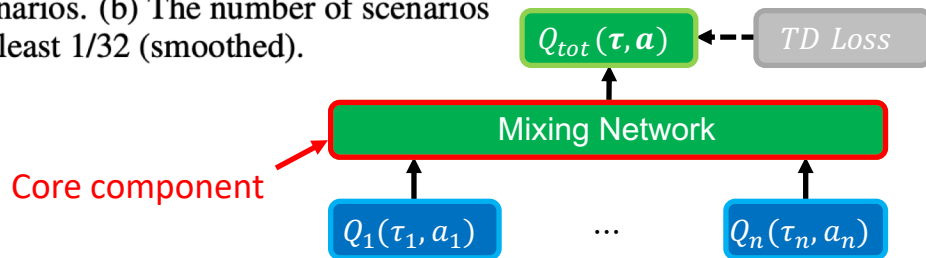


(a) Averaged test win rate



(b) # Maps best out of 17 scenarios

Figure 4: (a) The median test win %, averaged across all 17 scenarios. (b) The number of scenarios in which the algorithms' median test win % is the highest by at least 1/32 (smoothed).



# Thanks for your listening



Machine Intelligence Group



清华大学  
Tsinghua University

交叉信息研究院  
Institute for Interdisciplinary Information Sciences