# DFAC Framework: Factorizing the Value Function via Quantile Mixture for Multi-Agent Distributional Q-Learning

Wei-Fang Sun, Cheng-Kuang Lee, Chun-Yi Lee

# Outline

- Background & Motivation

- Proposed Method: DFAC

- Experiment Results: Outperform all baselines

# Outline

- **Background & Motivation**

- Proposed Method: DFAC
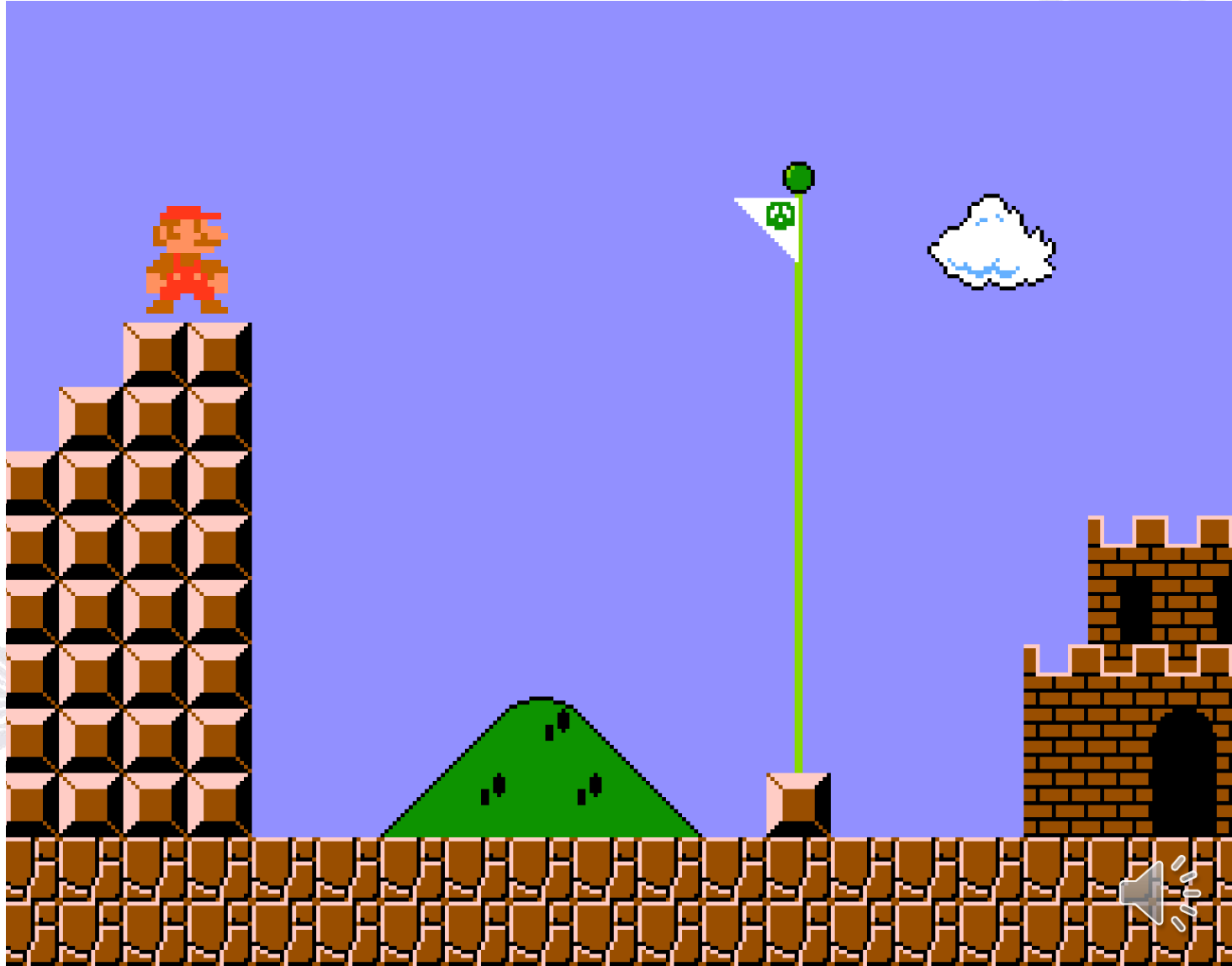
- Experiment Results: Outperform all baselines

Reinforcement Learning
(Focus on Q-Learning)

Reinforcement Learning
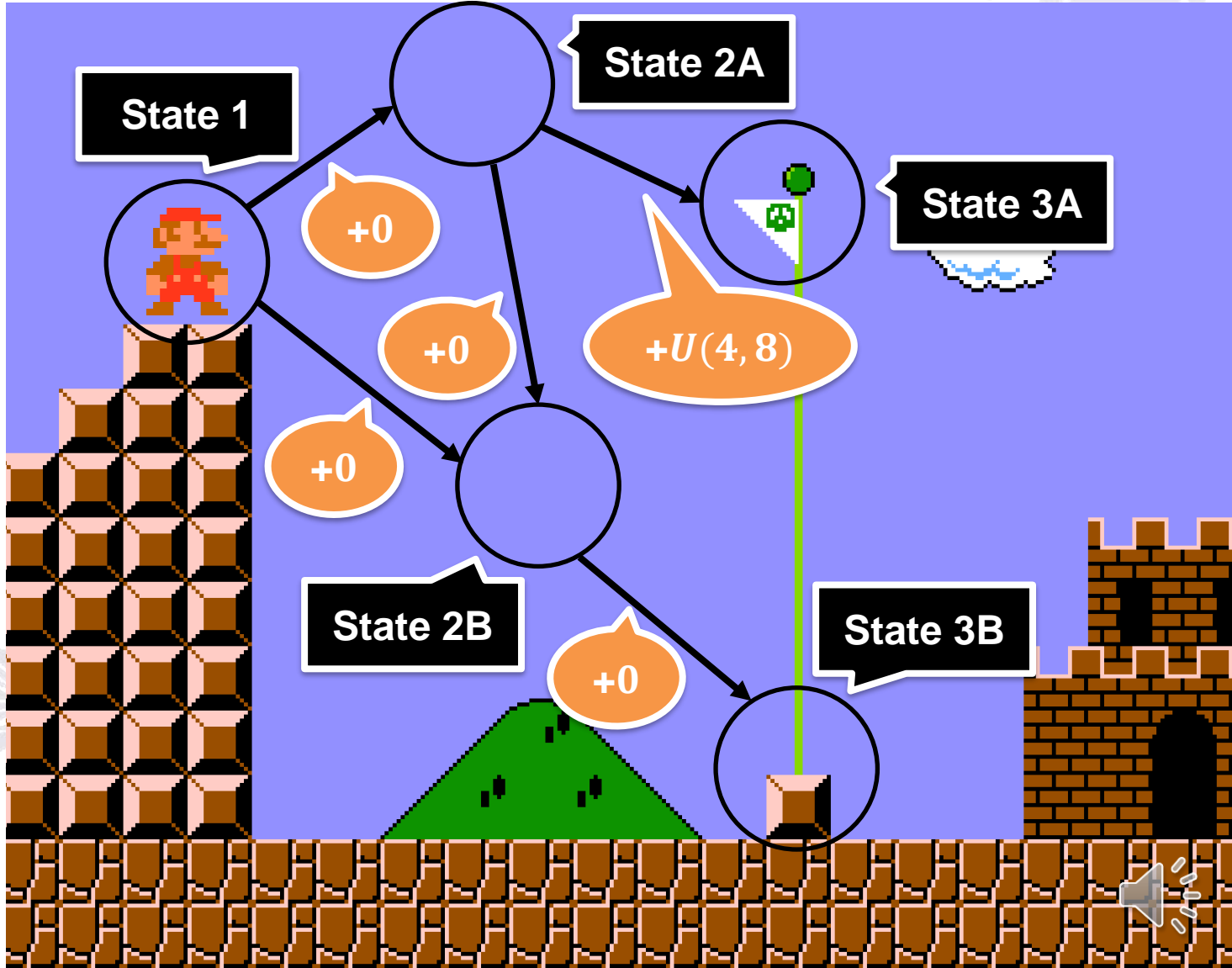(Focus on Q-Learning)

Single-Agent RL
(SARL)

Decentralized Partial Observable MDP

(Focus on Q-Learnin

IQN models CDF

Single-Agent RL
(SARL)

Distributional RL

Categorical Distribution (C51)

Implicit Quantile Network (IQN)

Enemy

Attack

Attack

Evade

Evade

Agent 1

Agent 2

Enemy attacking Agent

| | Attack 2 | |
|---|---|---|
| Attack 1 | +5 | +0 |
| Evade 1 | +10 | +5 |

Fully Cooperative
Team Reward (not Individual Reward)
⇒ Have Issues in Q-Learning

Decentralized Partial
Observable MDP

Fully Cooperative
Multi-Agent RL
(MARL)

Enemy

Attack

Attack

Evade

Evade

Agent 1

Agent 2

## Enemy attacking Agent 1 (Left)

| | Attack 2 | Evade 2 |
|---|---|---|
| Attack 1 | +5 | +0 |
| Evade 1 | +10 | +5 |

Fully Cooperative Multi-Agent RL (MARL)

## Enemy attacking Agent 2 (Right)

| | Attack 2 | Evade 2 |
|---|---|---|
| Attack 1 | +5 | +10 |
| Evade 1 | +0 | +5 |

Enemy

Agent 1

Agent 2

## Enemy attacking Agent 1 (Left)

|  | Attack 2 | Evade 2 |
|---|---|---|
| Attack 1 | +5 | +0 |
| Evade 1 | +10 | +5 |
| (Decomposed) | | |
|  | Attack 2 | Evade 2 |
| Attack 1 | +0 / +5 | +0 / +0 |
| Evade 1 | +5 / +5 | +5 / +0 |

## Enemy attacking Agent 2 (Right)

|  | Attack 2 | Evade 2 |
|---|---|---|
| Attack 1 | +5 | +10 |
| Evade 1 | +0 | +5 |
| (Decomposed) | | |
|  | Attack 2 | Evade 2 |
| Attack 1 | +5 / +0 | +5 / +5 |
| Evade 1 | +0 / +0 | +0 / +5 |

Decentralized Control → Value Factorization

Value Factorization

R... tion **u**:

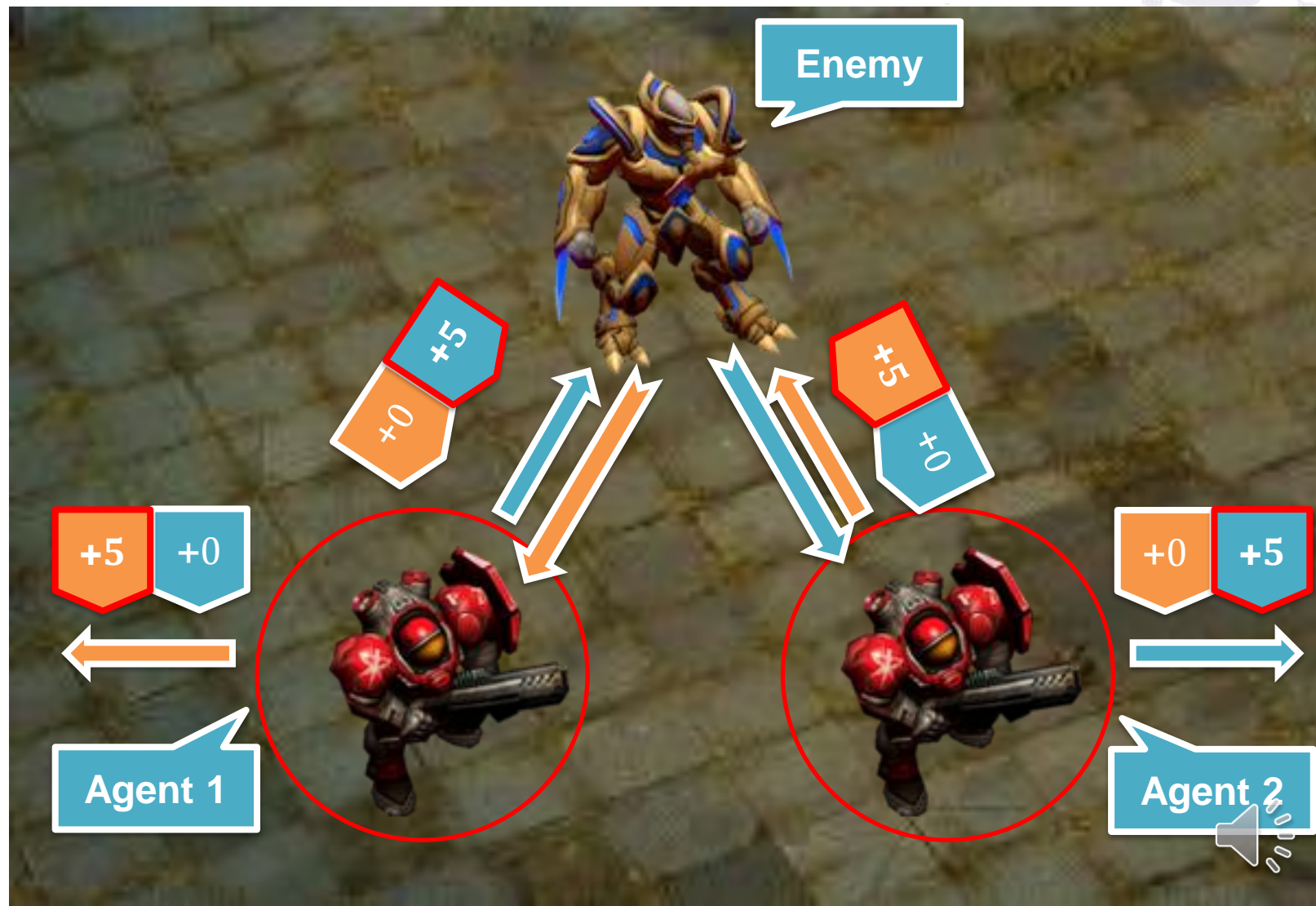$$Q_t(\mathbf{h}, \mathbf{u}) = \Psi(s, Q_1(h_1, u_1), \dots, Q_N(h_N, u_N))$$

| Enemy attacking Agent 1 (Left) | | |
|---|---|---|
| | Attack 2 | Evade 2 |
| Attack 1 | +5 | +0 |
| Evade 1 | **+10** | +5 |
| (Decomposed) | | |
| | Attack 2 | Evade 2 |
| Attack 1 | +0 / +5 | +0 / +0 |
| Evade 1 | **+5 / +5** | +5 / +0 |

| Enemy attacking Agent 2 (Right) | | |
|---|---|---|
| | Attack 2 | Evade 2 |
| Attack 1 | +5 | **+10** |
| Evade 1 | +0 | +5 |
| (Decomposed) | | |
| | Attack 2 | Evade 2 |
| Attack 1 | +5 / +0 | **+5 / +5** |
| Evade 1 | +0 / +0 | +0 / +5 |



Enemy

+5

+0

+5

+0

+5 +0

+0 +5

Agent 1

Agent 2

...noto
...letwo
...QMIX

# Outline

- Background & Motivation

- **Proposed Method: DFAC**

- Experiment Results: Outperform all baselines

# Theorem 1: Naïve generalization does not satisfy IGM

Given a factorization function $\Psi$ that satisfies IGM in the following form:

$$Q_{jt}(\mathbf{h}, \mathbf{u}) = \Psi(s, Q_1(h_1, u_1), \dots, Q_K(h_K, u_K))$$

The condition above is **not enough to guarantee** that

$$Z_{jt}(\mathbf{h}, \mathbf{u}) = \Psi(s, Z_1(h_1, u_1), \dots, Z_K(h_K, u_K))$$

Satisfies IGM for random variables.

# Theorem 2: Mean-Shape Decomposition satisfies IGM

Mean-Shape Decomposition:

$$Z_{\mathrm{jt}}(\mathbf{h},\mathbf{u}) = \mathbb{E}\big[Z_{\mathrm{jt}}(\mathbf{h},\mathbf{u})\big] + \big(Z_{\mathrm{jt}}(\mathbf{h},\mathbf{u}) - \mathbb{E}\big[Z_{\mathrm{jt}}(\mathbf{h},\mathbf{u})\big]\big)$$
$$= Z_{\mathrm{mean}}(\mathbf{h},\mathbf{u}) + Z_{\mathrm{shape}}(\mathbf{h},\mathbf{u})$$

where

- $Z_{\mathrm{mean}}(\mathbf{h},\mathbf{u}) = \Psi\big(s, Q_1(h_1, u_1), \ldots, Q_{\mathrm{K}}(h_{\mathrm{K}}, u_{\mathrm{K}})\big)$

- $Z_{\mathrm{shape}}(\mathbf{h},\mathbf{u}) = \Phi\big(s, Z_1(h_1, u_1), \ldots, Z_{\mathrm{K}}(h_{\mathrm{K}}, u_{\mathrm{K}})\big)$

- $\Psi$ satisfies IGM for $[Q_k]_{k=1}^{\mathrm{K}}$, $\mathrm{Var}(Z_{\mathrm{mean}}) = 0$, and $\mathbb{E}\big[Z_{\mathrm{shape}}\big] = 0$.

Mean-Shape decomposition is **guaranteed to satisfy IGM**.

# Theorem 3: Quantile Mixture have the form of sum of random variables
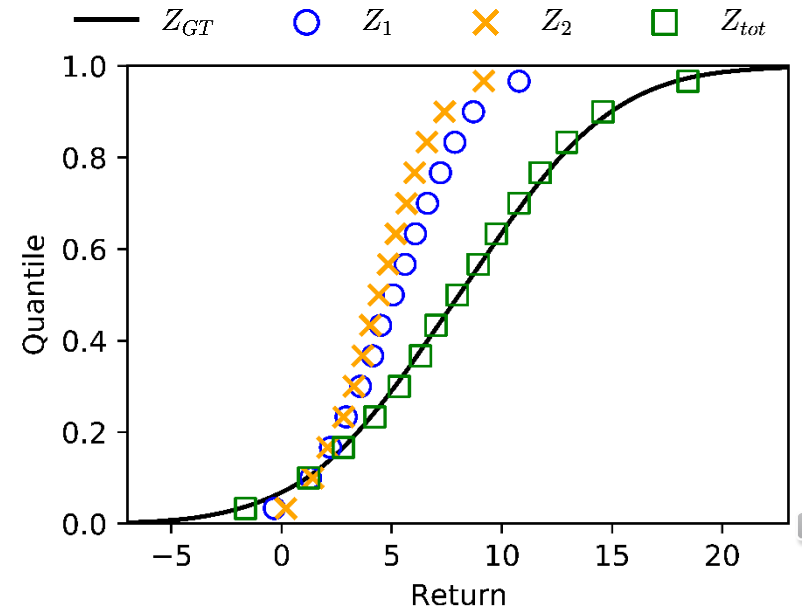
Given a Quantile Mixture:

$$F_Z^{-1} = \sum_{k=1}^{K} \beta_k \cdot F_{Z_k}^{-1}$$

where $\beta_k \geq 0, \forall k$. There exist corresponding $Z, [Z_k]_{k=1}^{K}$ that satisfies:

$$\textcolor{red}{Z = \sum_{k=1}^{K} \beta_k \cdot Z_k}$$

where the joint CDF of $[Z_k]_{k=1}^{K}$:



$$F_{\mathbf{Z}}(\mathbf{z}) = \min_{k} \left( F_{Z_k}^{-1}(z_k) \right)$$

# DFAC Framework

For approximating $Z_{\text{shape}} = \Phi$, we have the following choices:

(Assume we have $K$ agents and $N$ atoms/quantiles)

- C51 (models PMF)

  (convergence issue, not robust to hyperparameters, large network)
  - **Convolution** & Heuristic Projection ($\boldsymbol{O(KN^2)}$)
  - **FFT Convolution** + Heuristic Projection ($\boldsymbol{O(KN \log N)}$)

- IQN (models CDF)

  (convergence guarantee, robust to hyperparameters, light-weight)
  - **Quantile Mixture** ($\boldsymbol{O(KN)}$)

# Outline

- Background & Motivation

- Proposed Method: DFAC

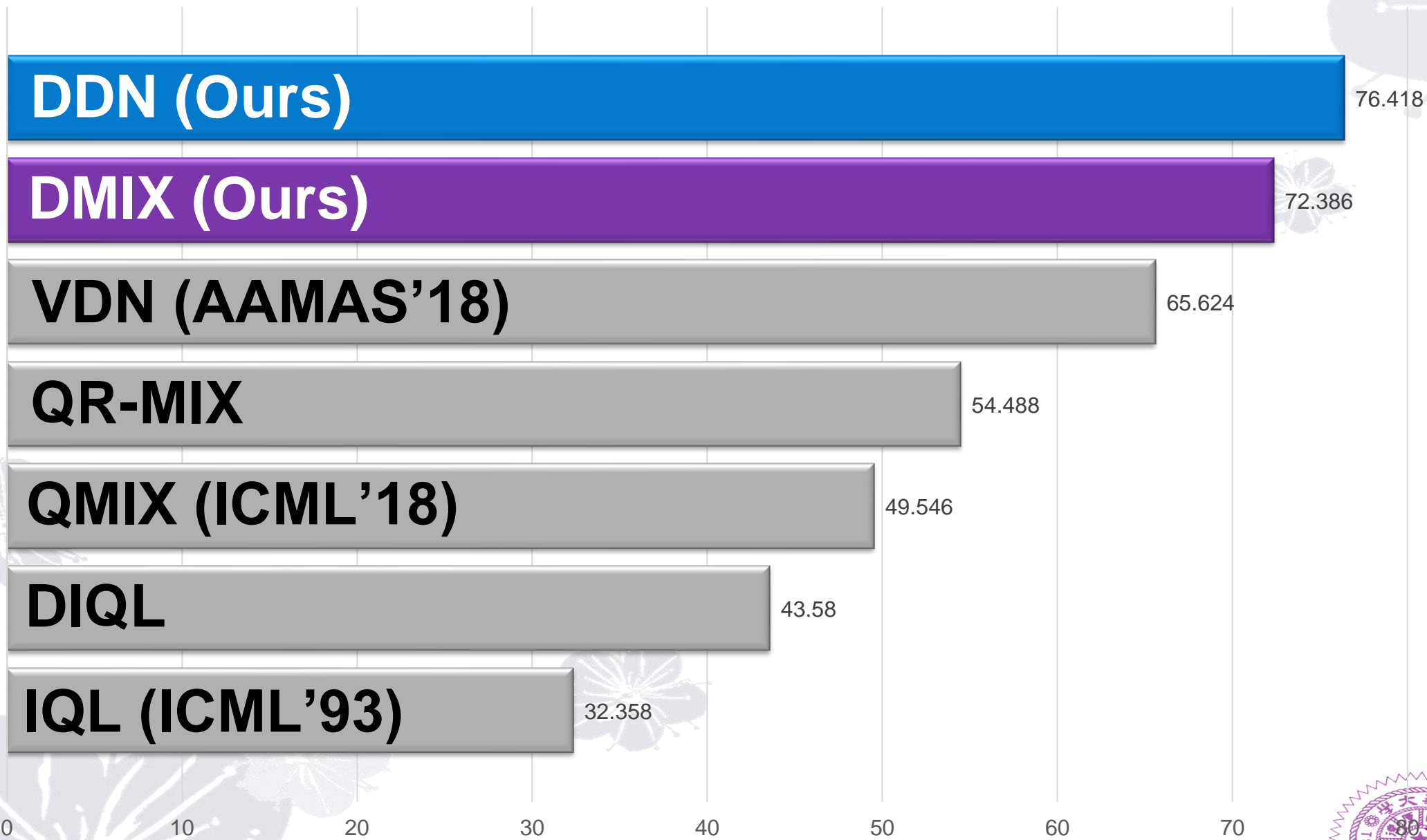- **Experiment Results: Outperform all baselines**

# Experiment Results

| | Map | IQL | VDN | QMIX | DIQL | **DDN** | **DMIX** |
|---|---|---|---|---|---|---|---|
| **Win Rate (%)** | 6h_vs_8z | 0.00% | 0.00% | 8.81% | 0.00% | **83.52%** | 68.75% |
| | 3s5z_vs_3s6z | 7.67% | 90.91% | 65.06% | 29.83% | **94.60%** | 90.62% |
| | MMM2 | 69.32% | 87.78% | 92.33% | 83.52% | **97.44%** | 95.17% |
| | 27m_vs_30m | 1.70% | 64.20% | 86.08% | 12.50% | **94.60%** | 86.08% |
| | corridor | 83.10% | 85.23% | 4.26% | 92.05% | **95.45%** | 90.06% |
| **Total Reward** | 6h_vs_8z | 13.96 | 15.49 | 14.02 | 14.98 | **19.32** | 17.81 |
| | 3s5z_vs_3s6z | 15.48 | 19.77 | 20.06 | 17.42 | 20.68 | **20.78** |
| | MMM2 | 17.47 | 19.32 | 19.45 | 19.21 | **21.06** | 19.69 |
| | 27m_vs_30m | 13.95 | 18.49 | 19.46 | 15.16 | **19.72** | 19.40 |
| | corridor | 19.30 | 19.38 | 13.44 | 19.57 | **19.97** | 19.61 |

# Win Rate (%)



| | Win Rate (%) |
|---|---|
| DDN (Ours) | 76.418 |
| DMIX (Ours) | 72.386 |
| VDN (AAMAS'18) | 65.624 |
| QR-MIX | 54.488 |
| QMIX (ICML'18) | 49.546 |
| DIQL | 43.58 |
| IQL (ICML'93) | 32.358 |

■DDN  ■DMIX  ▨VDN  ▨QR-MIX  ▨QMIX  ▨DIQL  ▨IQL

# Thank you!

For more information, please refer to the QR code below: