

# RiskQ : Risk-sensitive Multi-Agent Reinforcement Learning Value Factorization

Siqi Shen<sup>†</sup>, **Chennan Ma<sup>†</sup>**, Chao Li<sup>†</sup>, Weiquan Liu<sup>†</sup>,  
Yongquan Fu<sup>‡\*</sup>, Songzhu Mei<sup>‡</sup>, Xinwang Liu<sup>‡</sup>, Cheng Wang<sup>†</sup>

siqishen@xmu.edu.cn, chennanma@stu.xmu.edu.cn, chaoli@stu.xmu.edu.cn,  
yongquanf@nudt.edu.cn

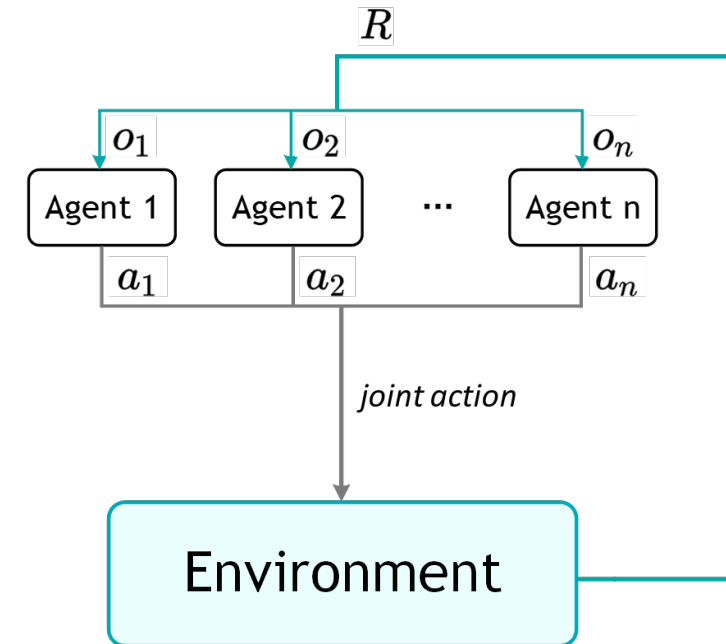
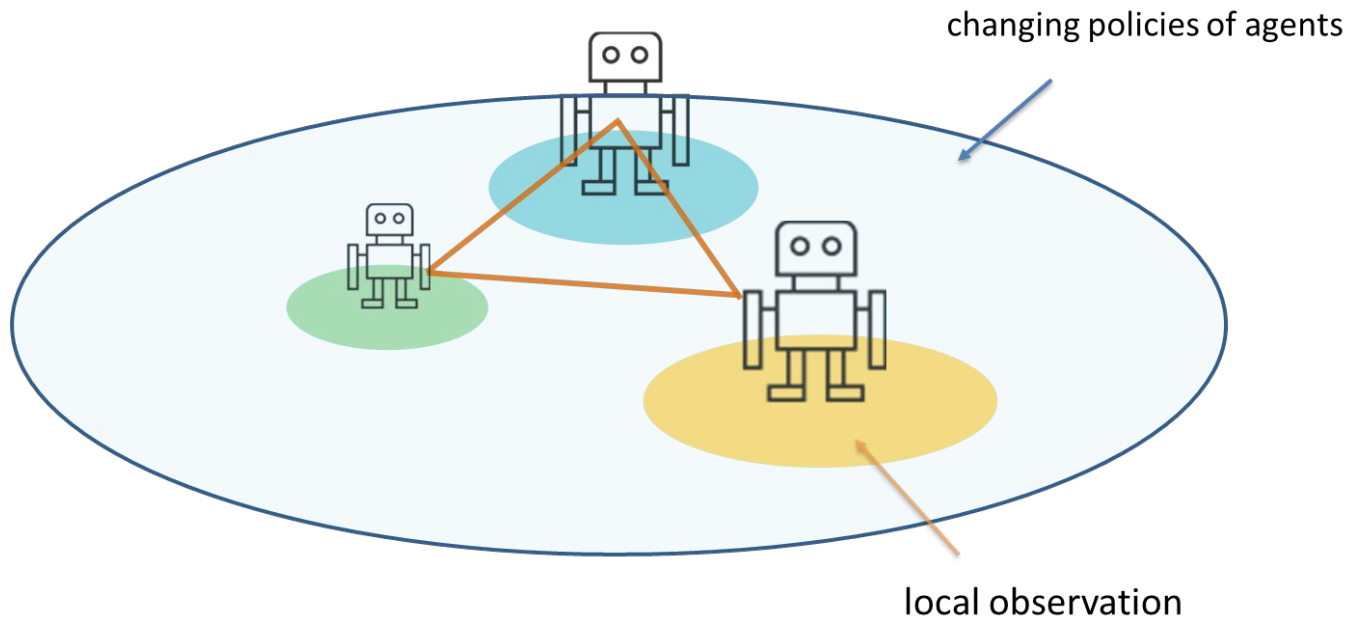
<sup>†</sup> Xiamen University

<sup>‡</sup> National University of Defense Technology

***NeurIPS 2023***

<https://arxiv.org/pdf/2311.01753.pdf>

# Challenges in MARL



Centralized Training with Decentralized Execution paradigm (CTDE)

# Value Factorization

- *Individual-Global-Max (IGM) principle*

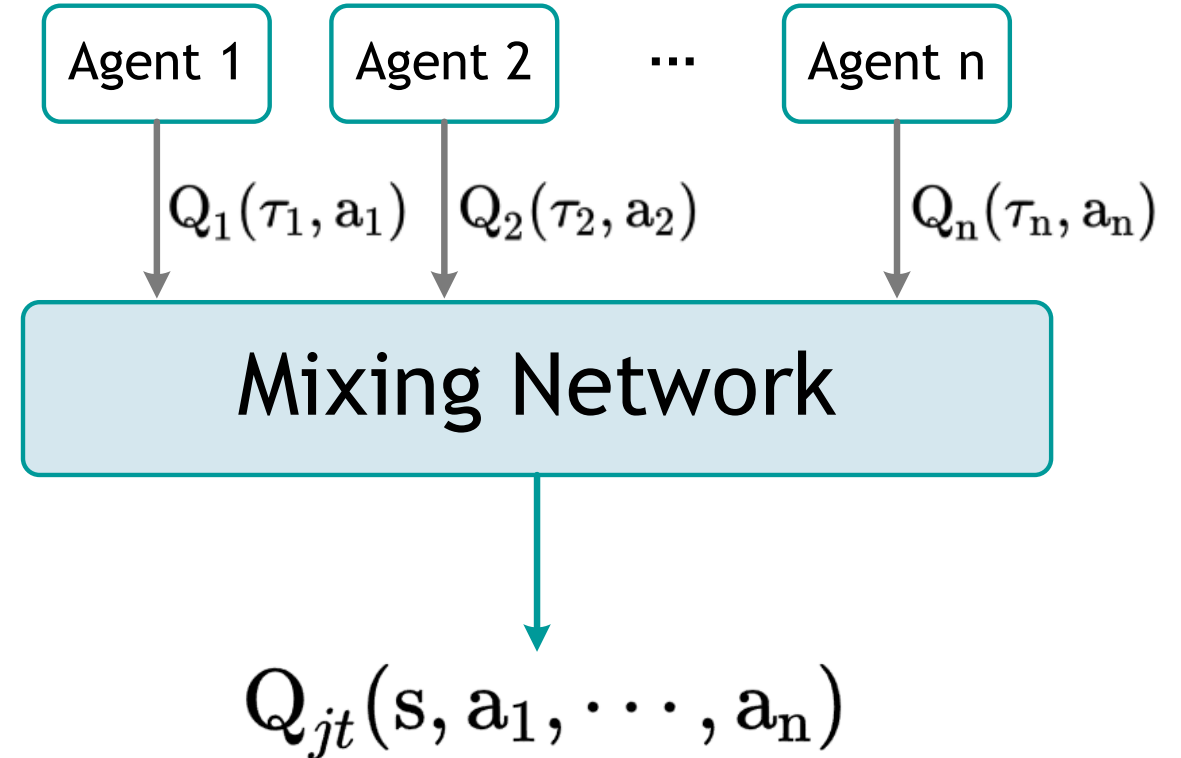
$$\arg \max_{\mathbf{u}} Q_{jt}(\boldsymbol{\tau}, \mathbf{u}) = \begin{pmatrix} \arg \max_{u_1} Q_1(\tau_1, u_1) \\ \vdots \\ \arg \max_{u_n} Q_n(\tau_n, u_n) \end{pmatrix}$$

**VDN**

$$Q_{jt}(\boldsymbol{\tau}, \mathbf{u}) = \sum_{i=1}^N Q_i(\tau_i, u_i)$$

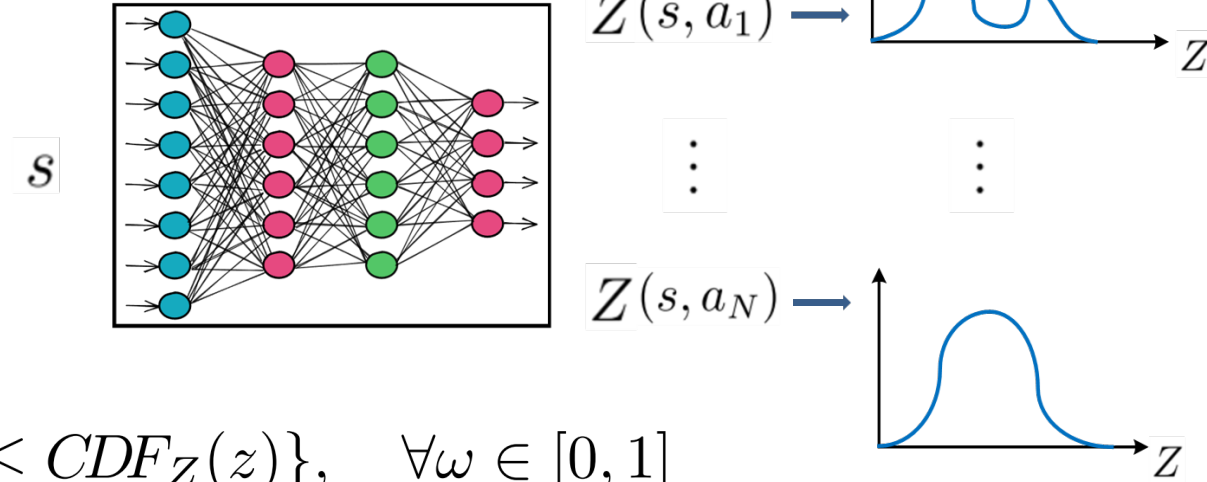
**QMIX**

$$\frac{\partial Q_{jt}(\boldsymbol{\tau}, \mathbf{u})}{\partial Q_i(\tau_i, u_i)} \geq 0, \quad \forall i \in \mathcal{N}$$



# Distributional RL

Value is actually a distribution



$$\theta_Z(\tau, \mathbf{u}, \omega) = \inf\{z \in \mathcal{R} : \omega \leq CDF_Z(z)\}, \quad \forall \omega \in [0, 1]$$

$$Z(\tau, \mathbf{u}) = \sum_{i=1}^n p_i(\tau, \mathbf{u}, \omega_i) \delta_{\theta(\tau, \mathbf{u}, \omega_i)}$$

- **Distributional IGM (DIGM) principle**

$$\begin{aligned} & \arg \max_{\mathbf{u}} \mathbb{E}[Z_{jt}(\tau, \mathbf{u})] \\ &= (\arg \max_{u_1} \mathbb{E}[Z_1(\tau_1, u_1)], \dots, \arg \max_{u_N} \mathbb{E}[Z_N(\tau_N, u_N)]) \end{aligned}$$

# Risk-sensitive RL

Risk-sensitive RL aims to optimize a **risk measure** based on a return distribution, rather than the expectation.

$$\Rightarrow \pi_{\psi_\alpha}(s) = \arg \max_u \psi_\alpha[Z(s, u)]$$

Risk measures:

- Value-at-risk (VaR)  $VaR_\alpha(Z(\tau, \mathbf{u})) = \theta(\tau, \mathbf{u}, \alpha)$

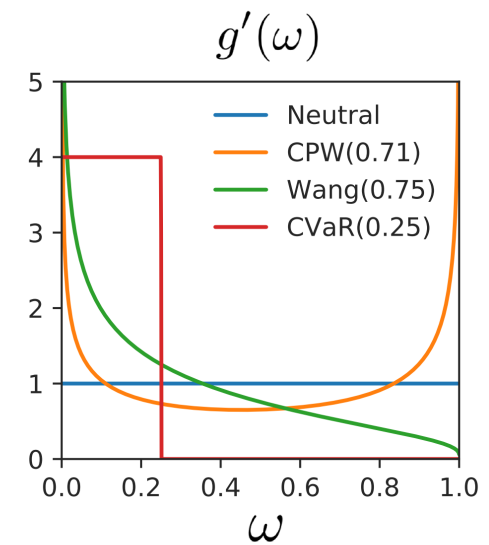
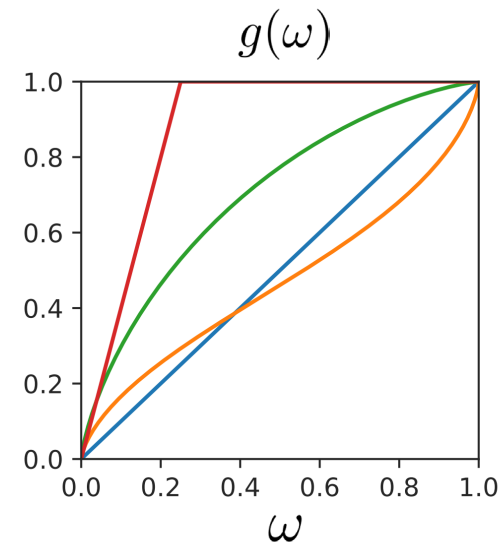
Distorted risk measure (DRM)  $\psi(Z) = \int_0^1 g'(\omega)\theta(\omega)d\omega$

- Conditional Value at Risk (CVaR)

$$CVaR_\alpha(Z) = \mathbb{E}_Z[z|z \leq \theta(\alpha)]$$

- Wang  $g(\omega) = \Phi(\Phi^{-1}(\omega) + \alpha)$

- CPW  $g(\omega) = \omega^\alpha / (\omega^\alpha + (1 - \omega)^\alpha)^{\frac{1}{\alpha}}$



# Motivation

➤ Risk-sensitive scenarios



Risk and Return



Risk-seeking



Risk-averse

- Most of the existing MARL value factorization methods do not extensively consider *risk*, which could impact their performance negatively in some **risk-sensitive scenarios**.
- How to **effectively** combine risk-sensitive reinforcement learning with MARL value factorization?

# RiskQ

## Risk-sensitive Individual-Global-Max (RIGM) Principle

**Definition 6 (RIGM).** Given a risk metric  $\psi_\alpha$ , a set of individual return distribution utilities  $[Z_i(\tau_i, u_i)]_{i=1}^N$ , and a joint state-action return distribution  $Z_{jt}(\boldsymbol{\tau}, \mathbf{u})$ , if the following conditions are satisfied:

$$\arg \max_{\mathbf{u}} \psi_\alpha[Z_{jt}(\boldsymbol{\tau}, \mathbf{u})] = (\arg \max_{u_1} \psi_\alpha[Z_1(\tau_1, u_1)], \dots, \arg \max_{u_N} \psi_\alpha[Z_N(\tau_N, u_N)]), \quad (7)$$

where  $\psi_\alpha : Z \times R \rightarrow R$  is a risk metric such as the VaR or a distorted risk measure,  $\alpha$  is its risk level. Then,  $[Z_i(\tau_i, u_i)]_{i=1}^N$  satisfy the RIGM principle with risk metric  $\psi_\alpha$  for  $Z_{jt}$  under  $\tau$ . We can state that  $Z_{jt}(\boldsymbol{\tau}, \mathbf{u})$  can be distributionally factorized by  $[Z_i(\tau_i, u_i)]_{i=1}^N$  with risk metric  $\psi_\alpha$ .

### The RIGM principle is a generalization of the DIGM and the IGM principle.

- $\psi = CVaR$  and  $\alpha = 1$ , RIGM principle  $\Rightarrow$  DIGM principle .
 
$$\arg \max_{\mathbf{u}} \mathbb{E}[Z_{jt}(\boldsymbol{\tau}, \mathbf{u})] = (\arg \max_{u_1} \mathbb{E}[Z_1(\tau_1, u_1)], \dots, \arg \max_{u_N} \mathbb{E}[Z_N(\tau_N, u_N)])$$
- If  $Z_i$  is a single Dirac Delta Distribution (value distribution  $Z_i$  becomes a single value, i.e.,  $Q_i$ ), and in this case ( $\psi = CVaR$  and  $\alpha = 1$ ), RIGM principle  $\Rightarrow$  IGM principle .
 
$$\arg \max_{\mathbf{u}} Q_{jt}(\boldsymbol{\tau}, \mathbf{u}) = \begin{pmatrix} \arg \max_{u_1} Q_1(\tau_1, u_1) \\ \vdots \\ \arg \max_{u_n} Q_n(\tau_n, u_n) \end{pmatrix}$$

# RiskQ Current value factorization methods can not satisfy RIGM principle

**Theorem 1.** *Given a deterministic joint action-value function  $Q_{jt}$ , a stochastic joint action-value function  $Z_{jt}$ , and a factorization function  $\Phi$  for deterministic utilities:*

$$Q_{jt}(\tau, u) = \Phi(Q_1(\tau_1, u_1), \dots, Q_n(\tau_n, u_n)) \quad (7)$$

*such that  $[Q_i]_{i=1}^n$  satisfy IGM for  $Q_{jt}$  under  $\tau$ , the following risk-sensitive distributional factorization:*

$$Z_{jt}(\tau, u) = \Phi(Z_1(\tau_1, u_1), \dots, Z_n(\tau_n, u_n)) \quad (8)$$

*is insufficient to guarantee that  $[Z_i]_{i=1}^n$  satisfy RIGM for  $Z_{jt}(\tau, u)$  with risk metric  $\psi_\alpha$ .*

**Theorem 2.** *Given a stochastic joint action-value function  $Z_{jt}$ , and a distributional factorization function  $\Phi$  for the stochastic utilities which satisfy the DIGM theorem, the following risk-sensitive distributional factorization:*

$$Z_{jt}(\tau, u) = \Phi(Z_1(\tau_1, u_1), \dots, Z_n(\tau_n, u_n)) \quad (9)$$

*is insufficient to guarantee that  $[Z_i]_{i=1}^n$  satisfy RIGM for  $Z_{jt}(\tau, u)$  with risk metric  $\text{VaR}_\alpha$ .*

**Theorem 3.** *DRIMA [14] does not guarantee adherence to the RIGM principle for CVaR metric.*

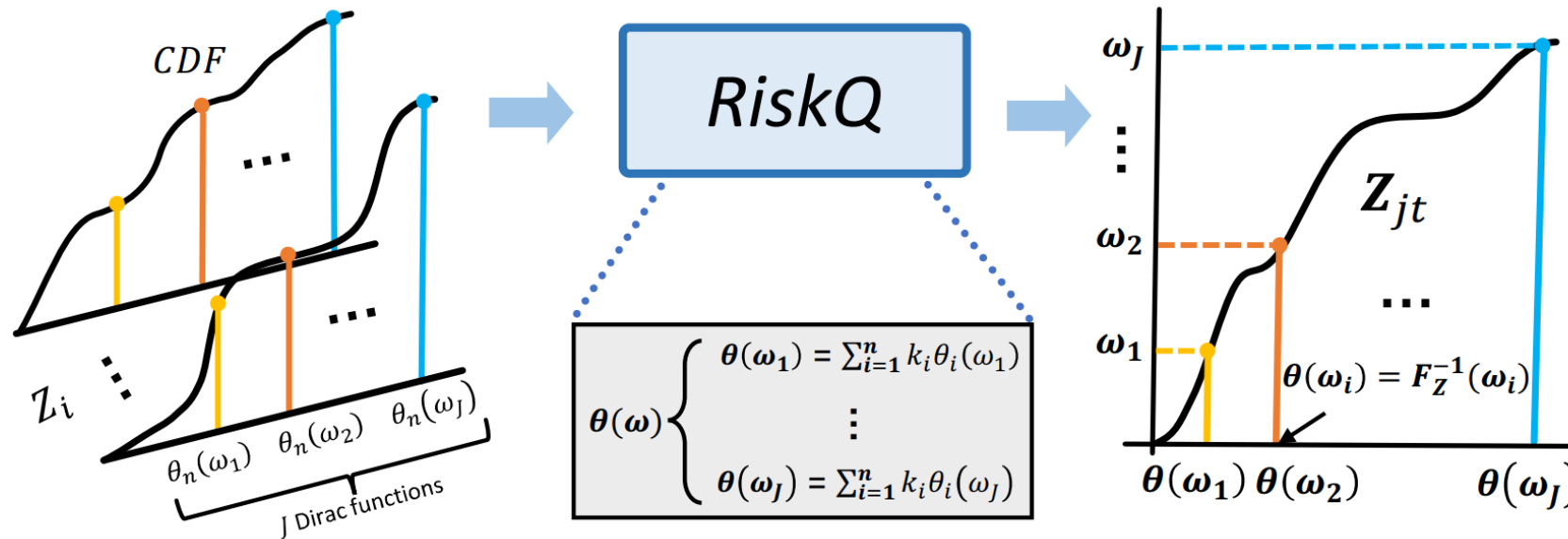


# RiskQ

## RiskQ satisfies RIGM principle

$$Z_{jt}(\boldsymbol{\tau}, \mathbf{u}) = \sum_{j=1}^J p_j(\boldsymbol{\tau}, \mathbf{u}, \omega_j) \delta_{\theta(\boldsymbol{\tau}, \mathbf{u}, \omega_j)}$$

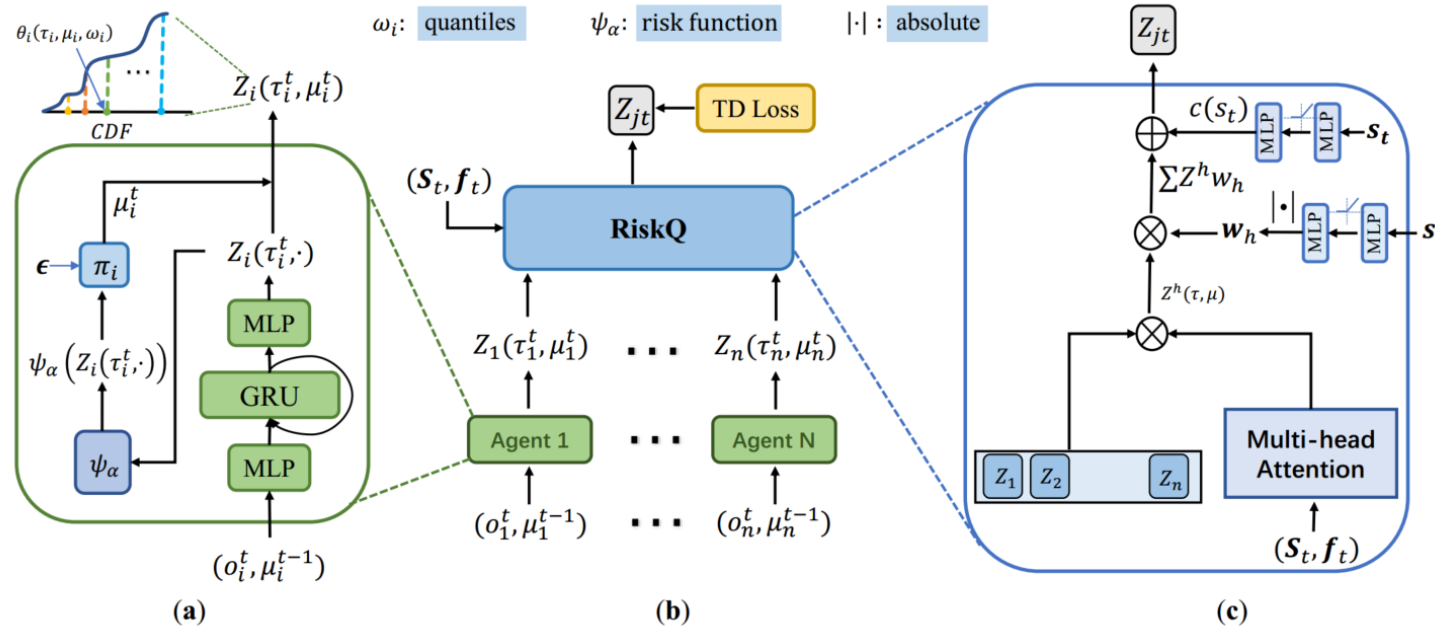
$$\theta(\boldsymbol{\tau}, \mathbf{u}, \omega_j) = \sum_{i=1}^N k_i \theta_i(\tau_i, u_i, \omega_j)$$



RiskQ overview: quantiles mixing for  $Z_{jt}$

# RiskQ

$$Z_{jt}(\boldsymbol{\tau}, \mathbf{u}) = \sum_{j=1}^J p_j(\boldsymbol{\tau}, \mathbf{u}, \omega_j) \delta_{\theta(\boldsymbol{\tau}, \mathbf{u}, \omega_j)} \quad \theta(\boldsymbol{\tau}, \mathbf{u}, \omega_j) = \sum_{i=1}^N k_i \theta_i(\tau_i, u_i, \omega_j)$$



target distribution:

$$y^k(\boldsymbol{\tau}^k, \mathbf{u}^k, \sigma) \triangleq r + \gamma Z_{jt}(\boldsymbol{\tau}^{k+1}, \tilde{\mathbf{u}}, \sigma^-)$$

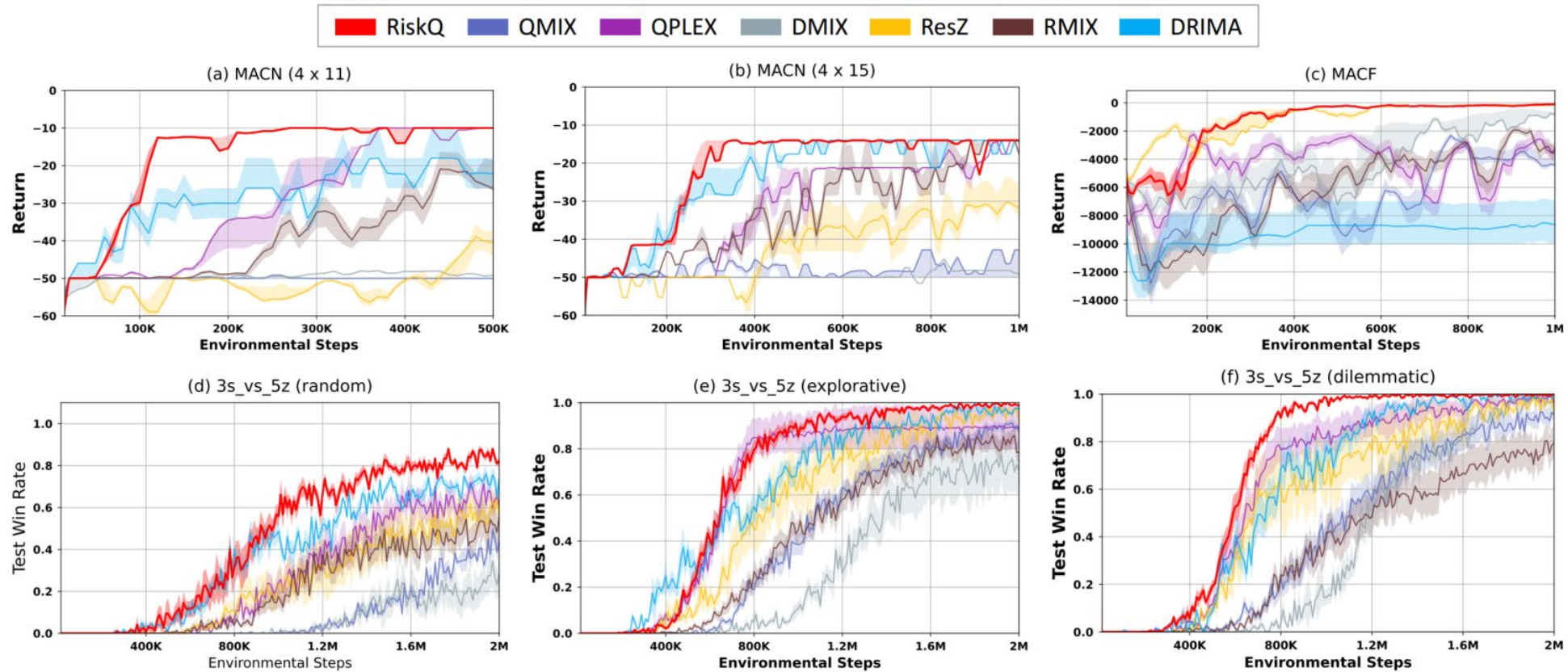
$$\tilde{\mathbf{u}} = [\tilde{u}_i]_{i=1}^N$$

$$\tilde{u}_i = \arg \max_{u_i} \psi_\alpha[Z_i(\tau_i^{k+1}, u_i)]$$

The framework of RiskQ

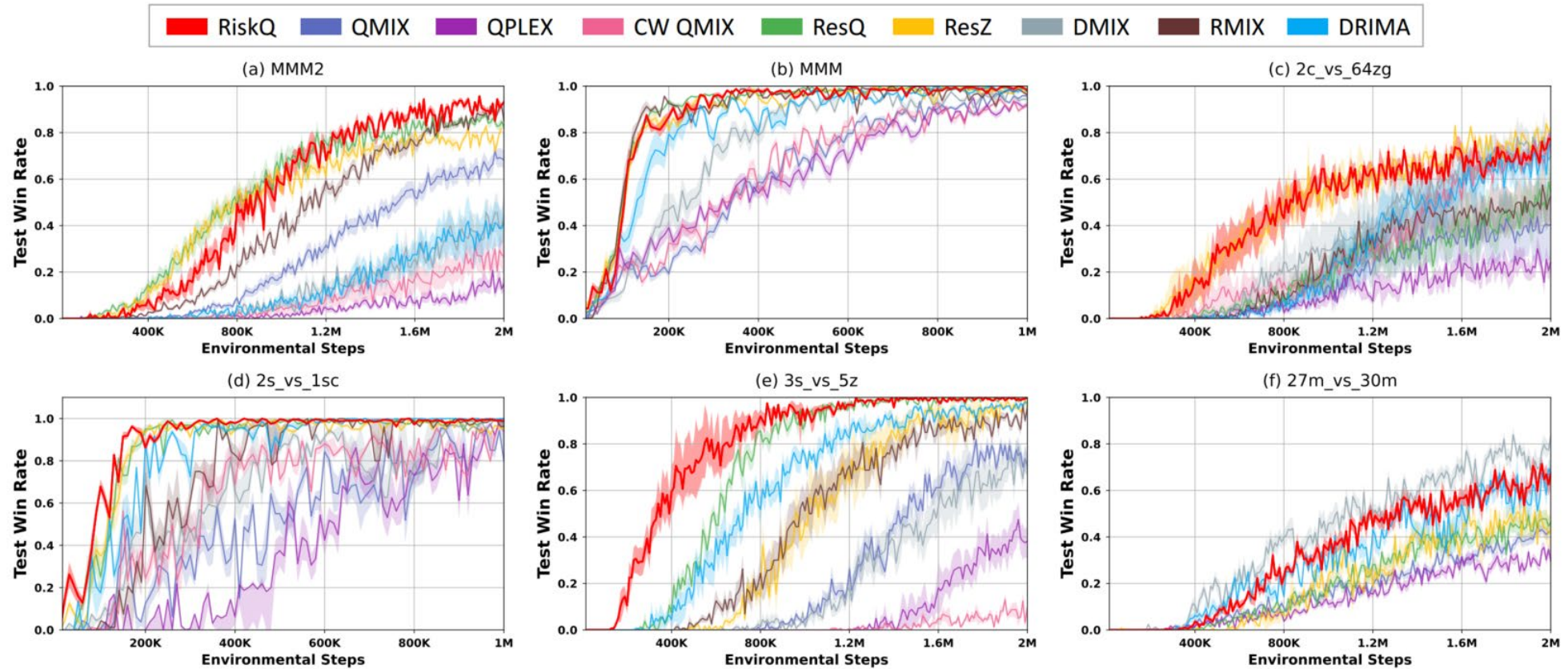
# Experiments — Risk-sensitive environments

RiskQ has better performance than other baselines in risk-sensitive settings



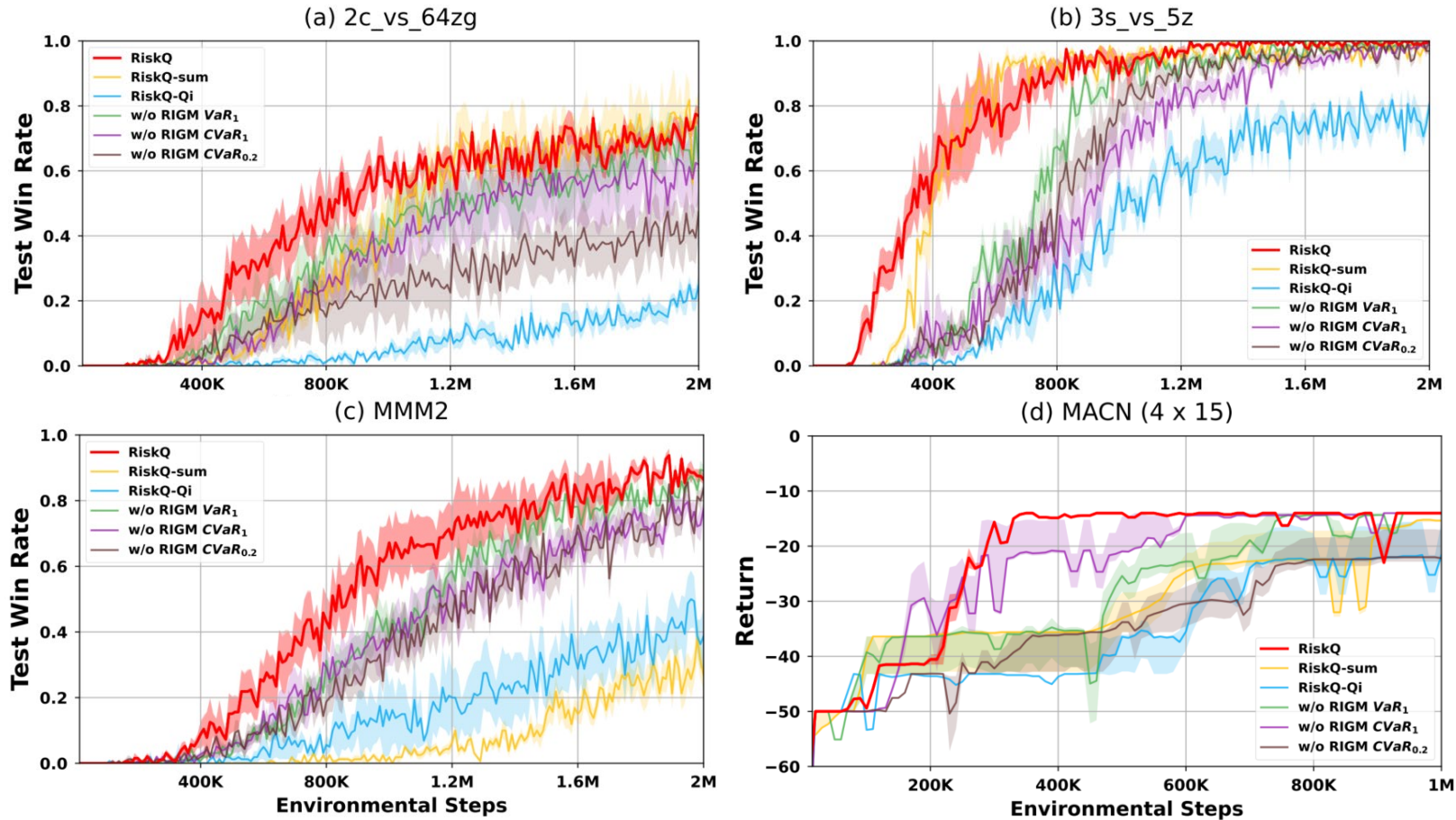
# Experiments — Starcraft II Multi-agent Challenge(SMAC)

RiskQ reaches the best win rate in most scenarios



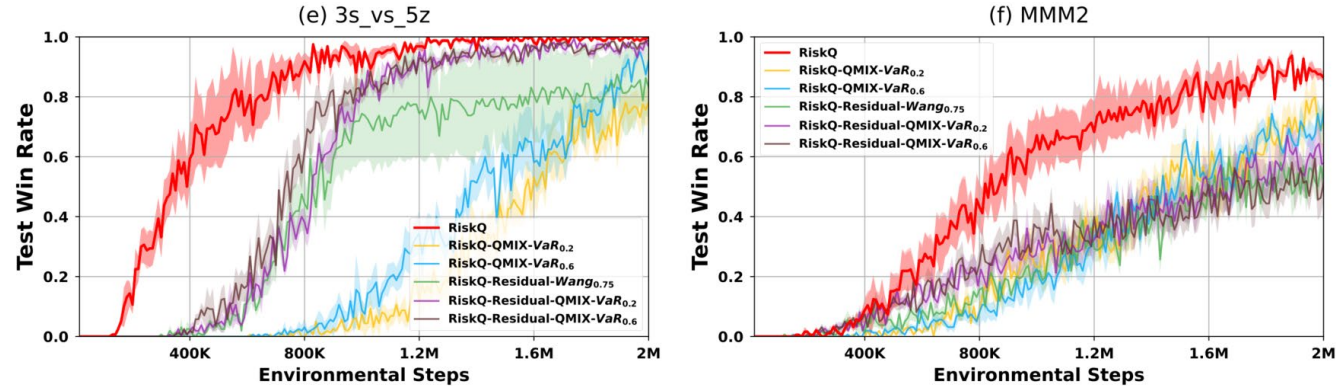
# Experiments — ablations

- It is important to satisfy the RIGM principle

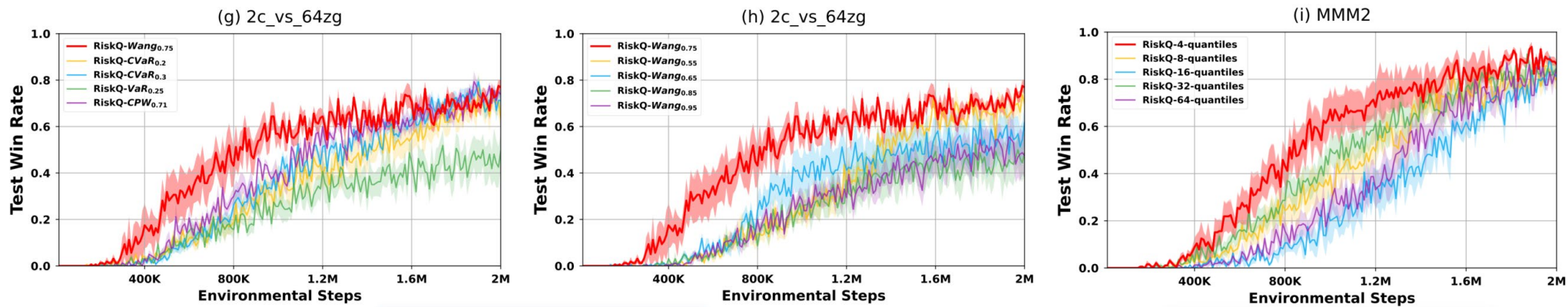


# Experiments — ablations

- The **representation limitations** of RiskQ do not significantly impact its performance



- Evaluate the impact of different **risk metrics, risk levels and number of percentiles**



# Summary

- RIGM principle, a generalization of IGM and DIGM principles.
- RiskQ, a value distribution factorization approach satisfying RIGM principle for Risk-sensitive Multi-Agent Reinforcement Learning problems
- Through extensive experiments, we show that RiskQ can obtain promising results.

*For more details, please check our project page:*

<https://github.com/xmu-rl-3dv/RiskQ>

**Contact us:**

[siqishen@xmu.edu.cn](mailto:siqishen@xmu.edu.cn)

[chennanma@stu.xmu.edu.cn](mailto:chennanma@stu.xmu.edu.cn)

[chaoli@stu.xmu.edu.cn](mailto:chaoli@stu.xmu.edu.cn)

[yongquanf@nudt.edu.cn](mailto:yongquanf@nudt.edu.cn)

