

The Capabilities and Limitations of Quantum Learning Models

Yuxuan Du

2021.04.10

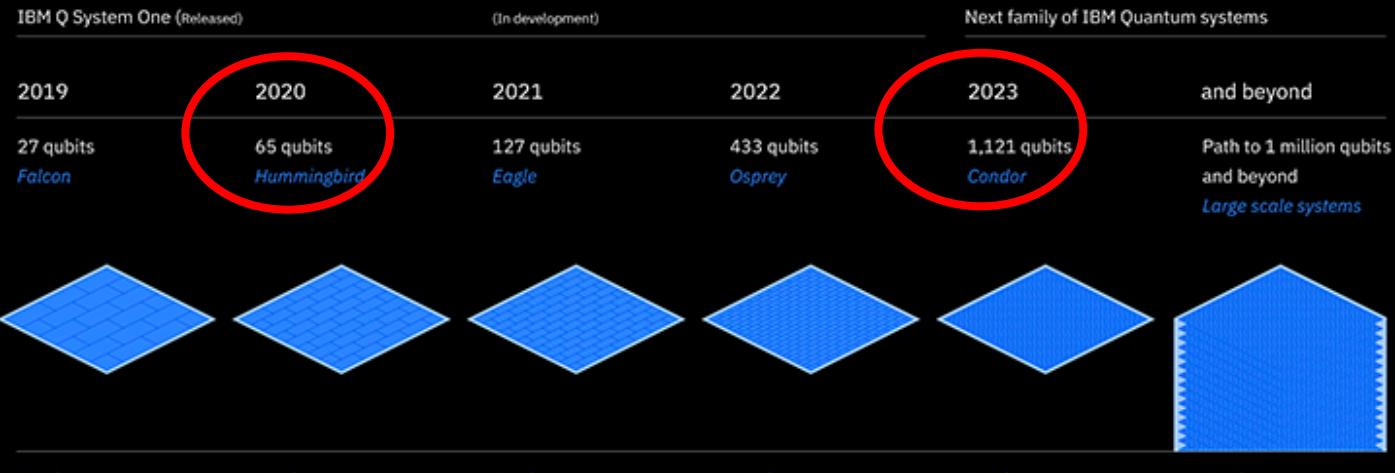
- I. The power of QNN in the view of optimization and learning theory [*On the learnability of quantum neural networks, arXiv:2007.12369*]
- II. QAS: An efficient scheme to enhance the trainability of QNN and suppress its error [*Quantum circuit architecture search: error mitigation and trainability enhancement for variational quantum solvers, arXiv:2010.10217*]
- III. The power of quantum kernels in the NISQ era [*Towards understanding the power of quantum kernels in the NISQ era, arXiv:2103.16774*]

Noisy intermediate-scale quantum (NISQ) era

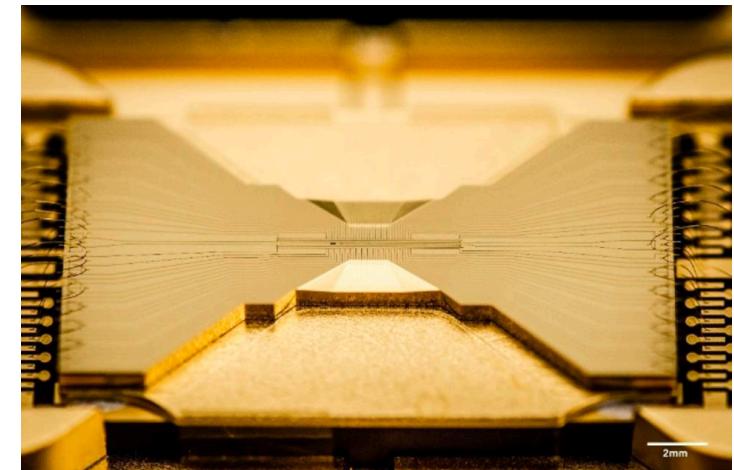
JDT 京东科技

Super-conducting

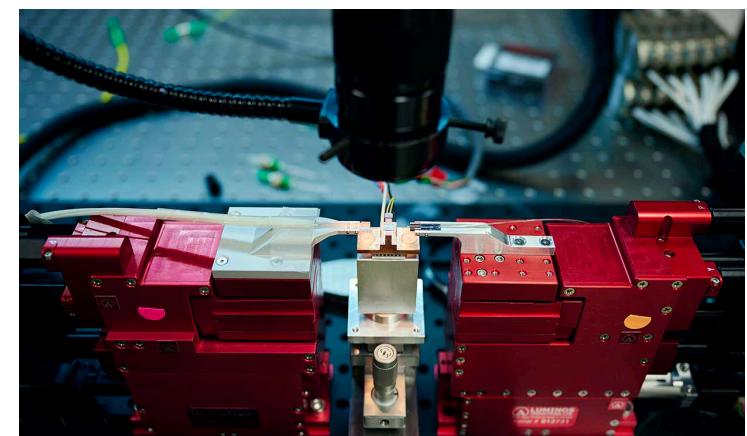
Scaling IBM Quantum technology



Trap-ion



Photonic chips



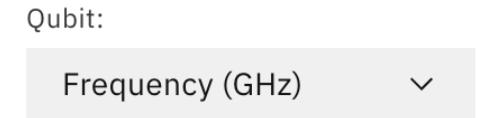
<https://www.ibm.com/blogs/research/2020/09/ibm-quantum-roadmap/>, <https://www.forbes.com/sites/ibm/2020/01/16/the-quantum-computing-era-is-here-why-it-matters-and-how-it-may-change-our-world/>, <https://ai.googleblog.com/2019/10/quantum-supremacy-using-programmable.html>, <https://www.forbes.com/sites/moorinsights/2020/10/07/ionq-releases-a-new-32-qubit-trapped-ion-quantum-computer-with-massive-quantum-volume-claims>, <https://spectrum.ieee.org/tech-talk/computing/hardware/photonic-quantum>

Noisy intermediate-scale quantum (NISQ) era

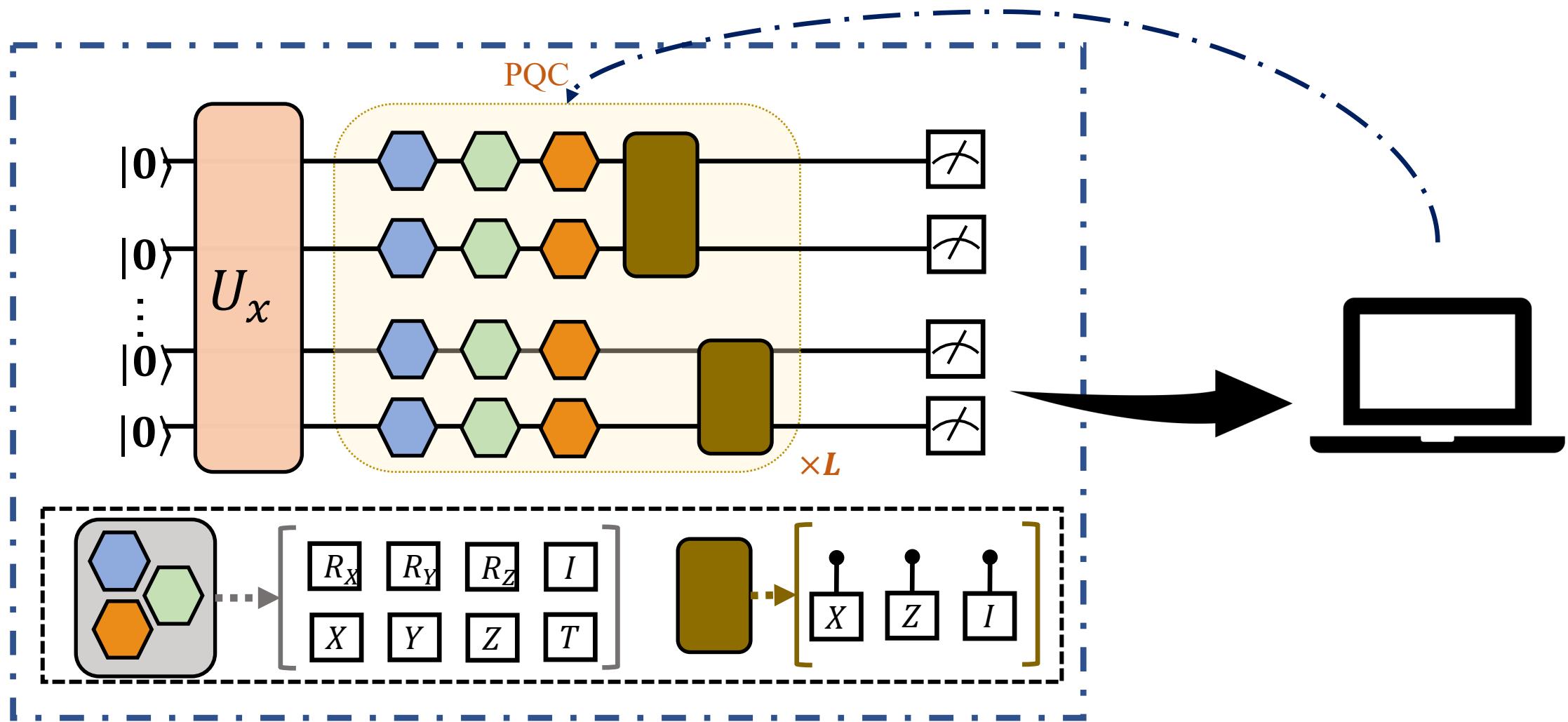
JDT 京东科技

The **caveats** of NISQ processors are:

- a limited number of qubits;
- shallow circuit depth;
- connectivity restriction;
- unavoidable system noise.

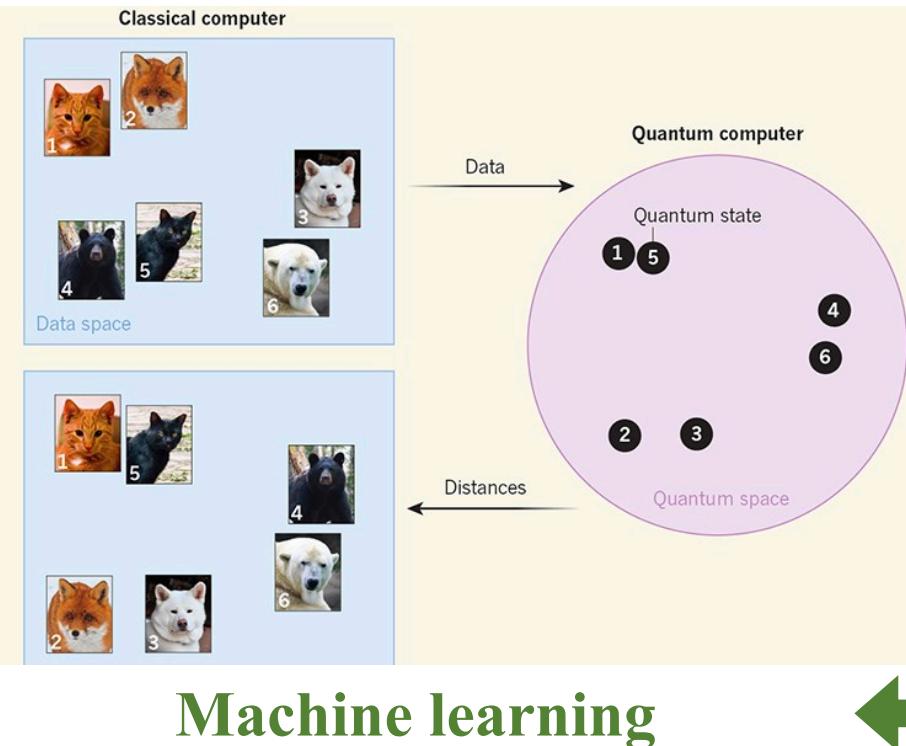


Variational quantum algorithm

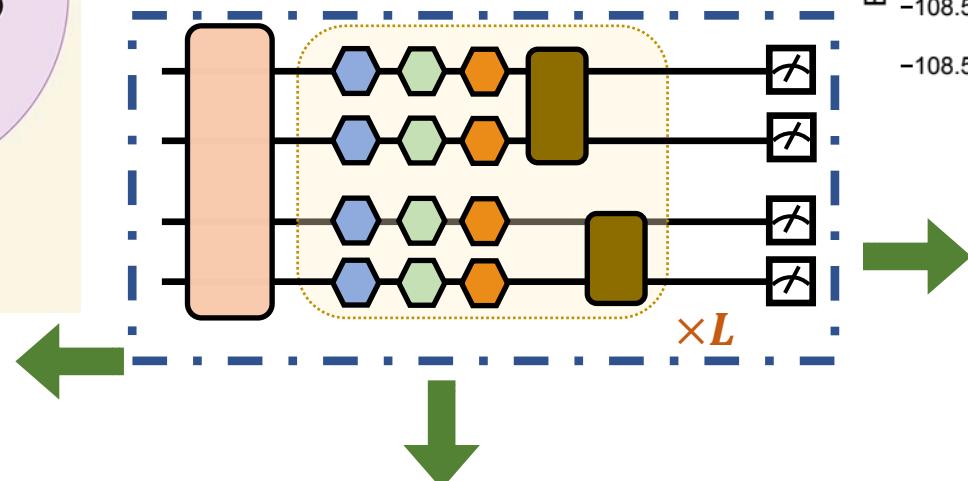


Noisy intermediate-scale quantum (NISQ) applications

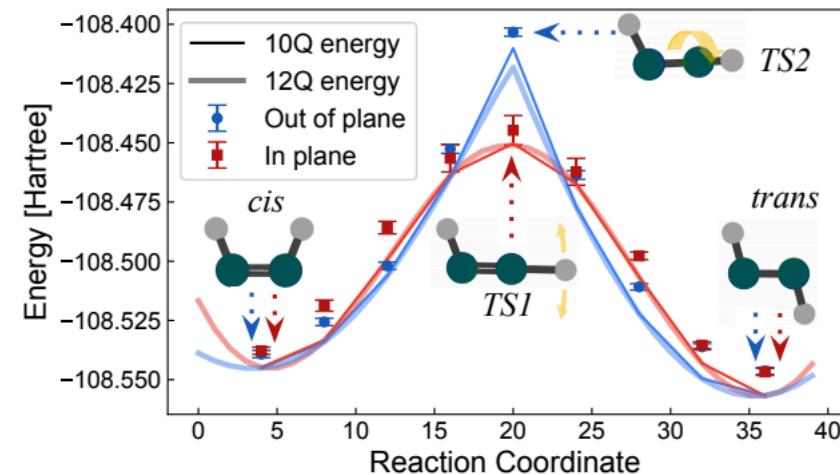
JDT 京东科技



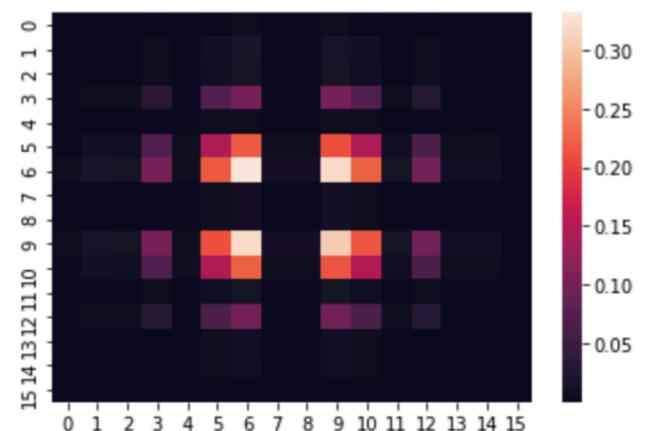
Variational quantum algorithm



Quantum information processing,
Quantum metrology, ...



Chemistry



Problem setup

ERM: The empirical risk minimization (ERM) principle is a learning paradigm that has been broadly employed to benchmark the performance of **supervised learning** algorithms.

When QNN is employed to tackle ERM, the optimization yields:

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta} \in \mathcal{C}} \mathcal{L}(\boldsymbol{\theta}, \mathbf{z}) := \frac{1}{n} \sum_{j=1}^n l(y_j, \hat{y}_j) + r(\boldsymbol{\theta}).$$

$\mathbf{z} = \{\mathbf{z}_j\}_{j=1}^N \in \mathcal{Z}$: Input dataset with N samples;

$\mathbf{z}_j = (\mathbf{x}_j, y_j)$: The j-th sample, the feature vector as $\mathbf{x}_j \in \mathbb{R}^{D_c}$, the label as $y_j \in \mathbb{R}^1$,

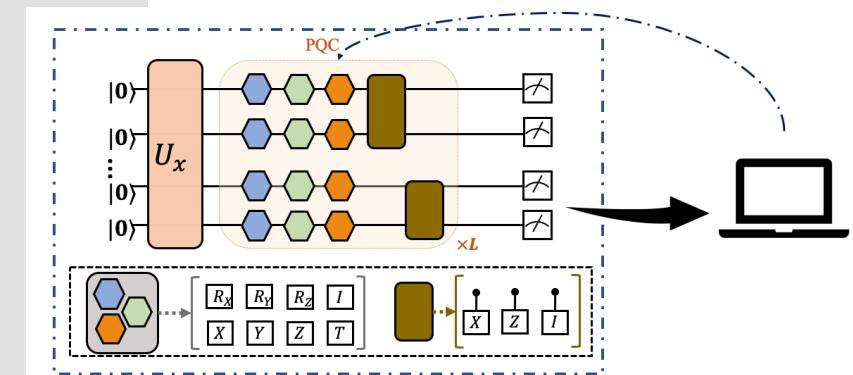
$\hat{y}_j = \text{Tr}((\Pi|\psi(\boldsymbol{\theta})\rangle\langle\psi(\boldsymbol{\theta})|))$: The predicted label of the j-th sample;

$|\psi(\boldsymbol{\theta})\rangle$: quantum states prepared by $U(\boldsymbol{\theta})$ and \mathbf{z}_j ;

$U(\boldsymbol{\theta})$: The trainable quantum circuit as.

$l(\cdot, \cdot)$: the loss function measures the difference between y_j and \hat{y}_j ;

$r(\boldsymbol{\theta})$: the regularizer.

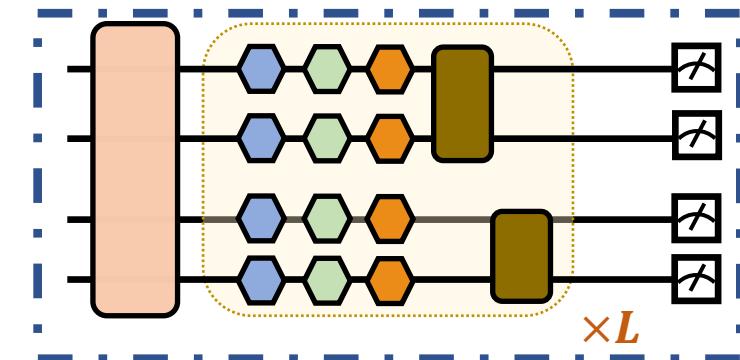


Challenges of QNN

$$\theta^* = \arg \min_{\theta \in \mathcal{C}} \mathcal{L}(\theta, z) := \frac{1}{n} \sum_{j=1}^n l(y_j, \hat{y}_j) + r(\theta).$$

The **difficulties** to theoretically understand QNN are:

- the **versatile structures**;
- the **non-convex optimization landscapes**;
- the **unavoidable gate noise** and **measurement errors**.



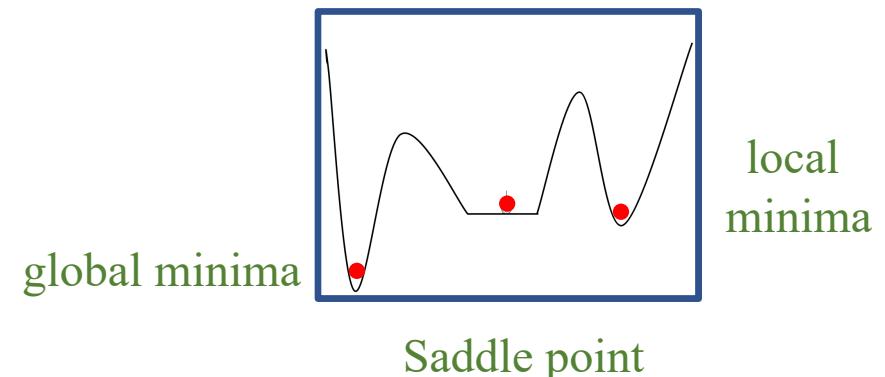
Part I. Understanding the power of QNN

Challenges of QNN

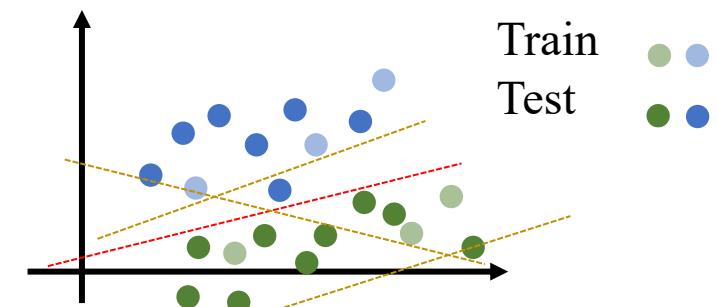
Key issue: what is the **learnability** of QNN, where

$$\text{learnability} = \text{'trainability'} + \text{'generalization'}$$

- *Trainability* \leftrightarrow stationary point;

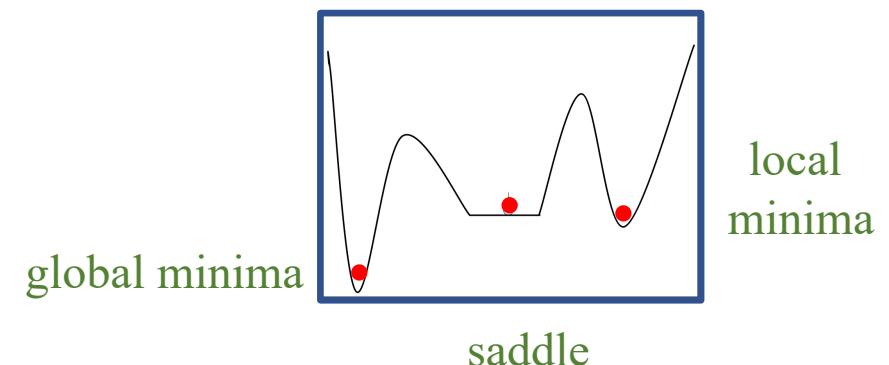


- *Generalization* \leftrightarrow find a good hypothesis in a poly sample complexity.



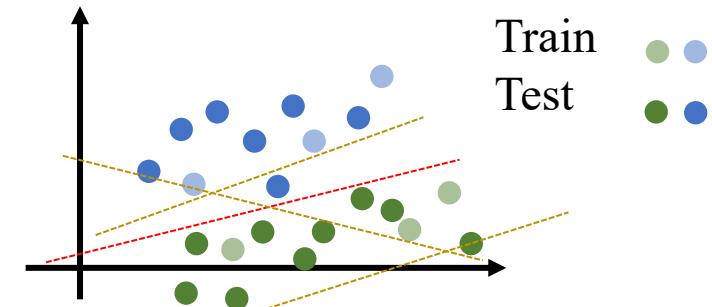
Trainability:

- ❖ **Evaluation** of the performance of various QNN-based learning algorithms;
- ❖ **Guide** to design better quantum learning protocols to avoid barren plateaus.



Generalization:

- ❖ Answer whether there exists any class of concepts that can be **efficiently** learned by (noisy) QNN but are **computationally hard** for the classical learning models.
- ❖ Use QNN implemented on NISQ devices to accomplish certain tasks with **theoretical advantages**.



Problem setup of the **trainability** of QNN towards ERM

JDT 京东科技

QNN towards ERM:

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta} \in \mathcal{C}} \mathcal{L}(\boldsymbol{\theta}, \mathbf{z}) := \frac{1}{n} \sum_{j=1}^n l(y_j, \hat{y}_j) + r(\boldsymbol{\theta}).$$

Setting: the mean square error loss with $l(y_j, \hat{y}_j) = (y_j - \hat{y}_j)^2$;

the L2-norm regularizer $r(\boldsymbol{\theta}) = \frac{\lambda \|\boldsymbol{\theta}\|^2}{2}$,

depolarization noise.

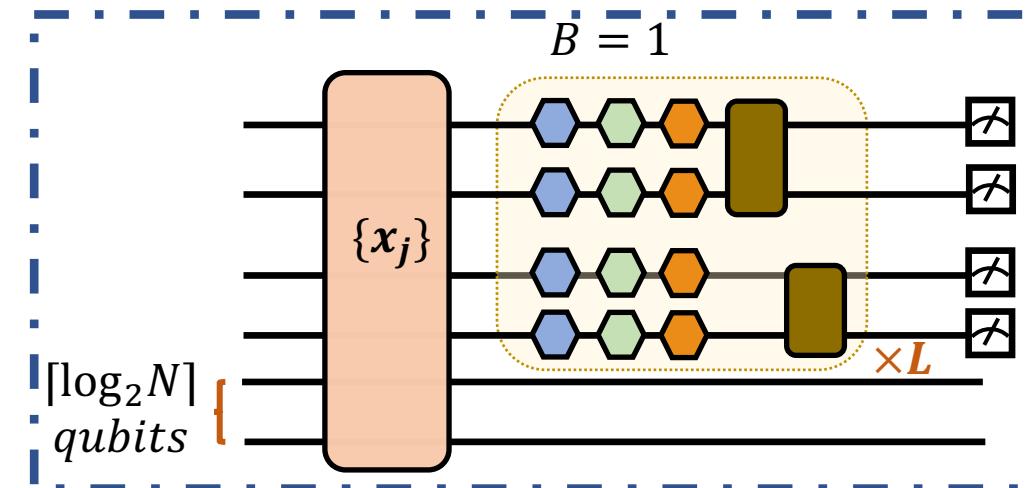
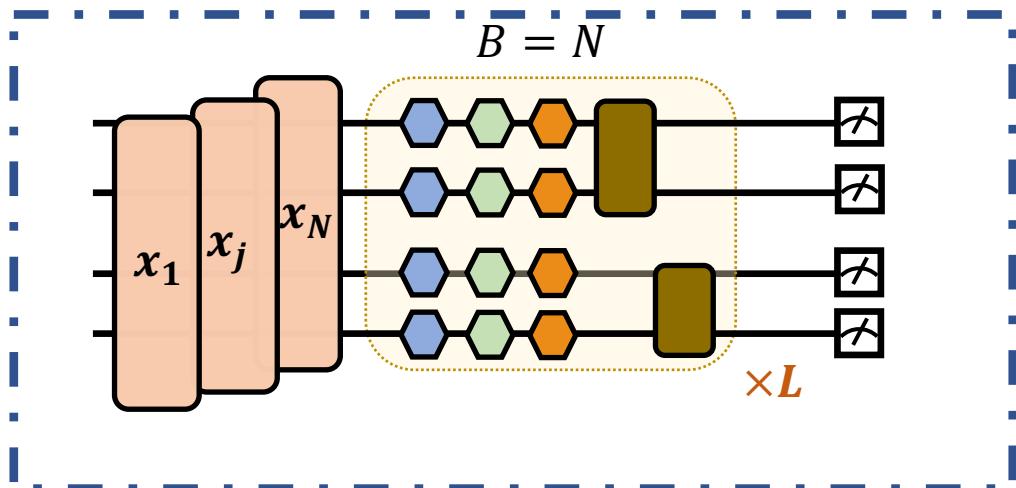
Our analysis can be easily generalized to other loss functions, regularizer, noise model.

Problem setup of the trainability of QNN towards ERM

The optimization rule is the batch gradient descent method. At the t -th iteration, the parameters are updated to

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \frac{\eta}{B} \sum_{i=1}^B \nabla \mathcal{L}(\boldsymbol{\theta}, B_i),$$

where η is the learning rate and $\nabla \mathcal{L}(\boldsymbol{\theta}, B_i)$ is obtained by *the parameter shift rule*.



Problem setup of the trainability of QNN towards ERM

JDT 京东科技

The trainability (convergence performance) of QNN is measured by the following two standard utility metrics:

$$R_1(\boldsymbol{\theta}^{(T)}) := \mathbb{E} \left[\left\| \nabla \mathcal{L}(\boldsymbol{\theta}^{(T)}) \right\|^2 \right],$$

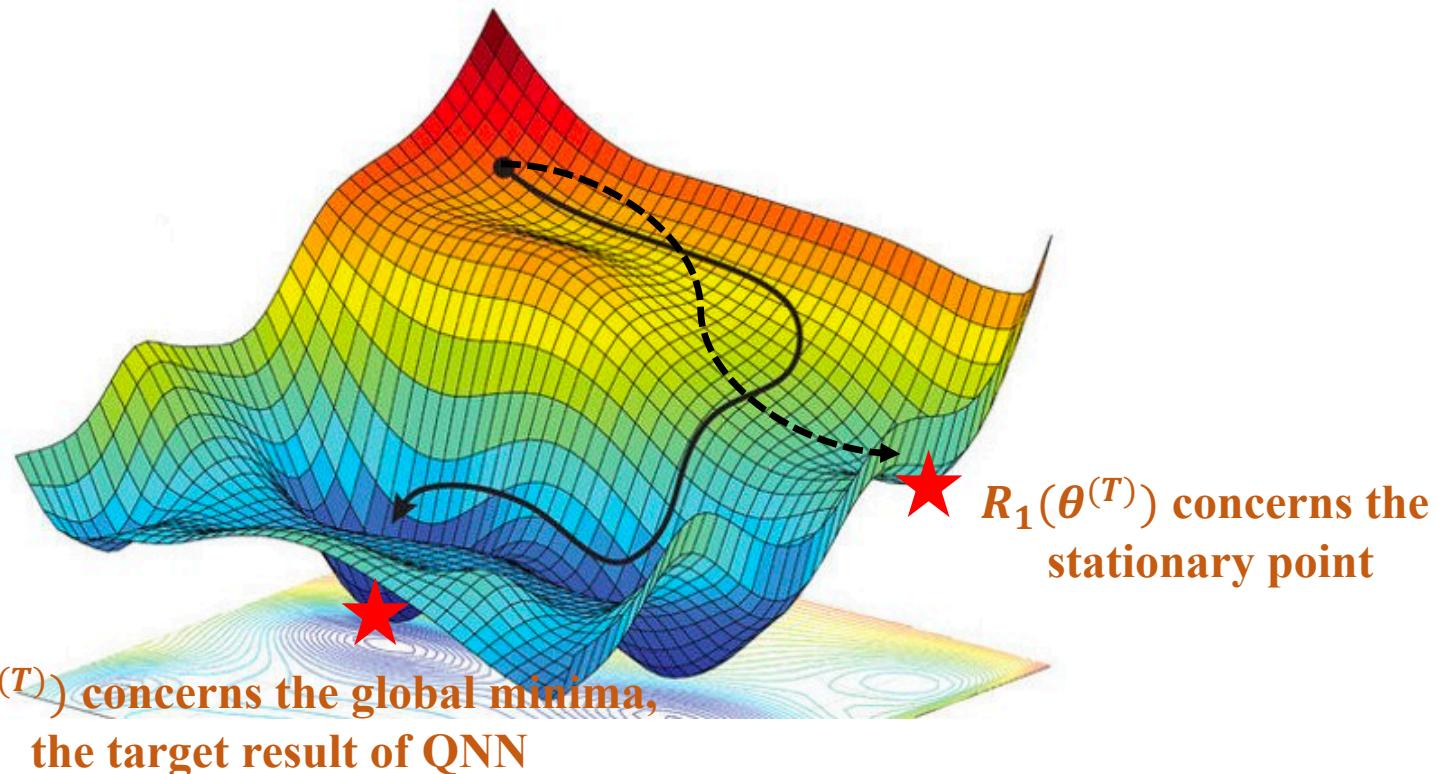
$$R_2(\boldsymbol{\theta}^{(T)}) := \mathbb{E}[\mathcal{L}(\boldsymbol{\theta}^{(T)})] - \mathcal{L}(\boldsymbol{\theta}^*),$$

where the expectation is taken over the randomness of QNN resulted from the measurement error and gate noise.

Problem setup of the trainability of QNN towards ERM

JDT 京东科技

- $R_1(\boldsymbol{\theta}^{(T)})$ evaluates how far QNN is away from **the stationary point**, i.e., $\|\nabla \mathcal{L}(\boldsymbol{\theta}^{(T)}, \mathbf{z})\| = 0$.
- $R_2(\boldsymbol{\theta}^{(T)})$ evaluates how far QNN is away from **the optimal result $\mathcal{L}(\boldsymbol{\theta}^*, \mathbf{z})$** .



Theorem 1. Let K be the number of measurements, L be the circuit depth, p be the depolarization noise of quantum gates, and B be the batch size. QNN outputs $\boldsymbol{\theta}^{(T)} \in \mathbb{R}^d$ after training T iterations with the utility bound

$$R_1(\boldsymbol{\theta}^{(T)}) \leq \tilde{O}\left(\text{poly}\left(\frac{d}{TBK(1-p)^L}\right)\right).$$

When λ satisfies a technical assumption, QNN outputs $\boldsymbol{\theta}^{(T)} \in \mathbb{R}^d$ after $T = \tilde{O}(d^3/(1-p)^L)$ training iterations with the utility bound

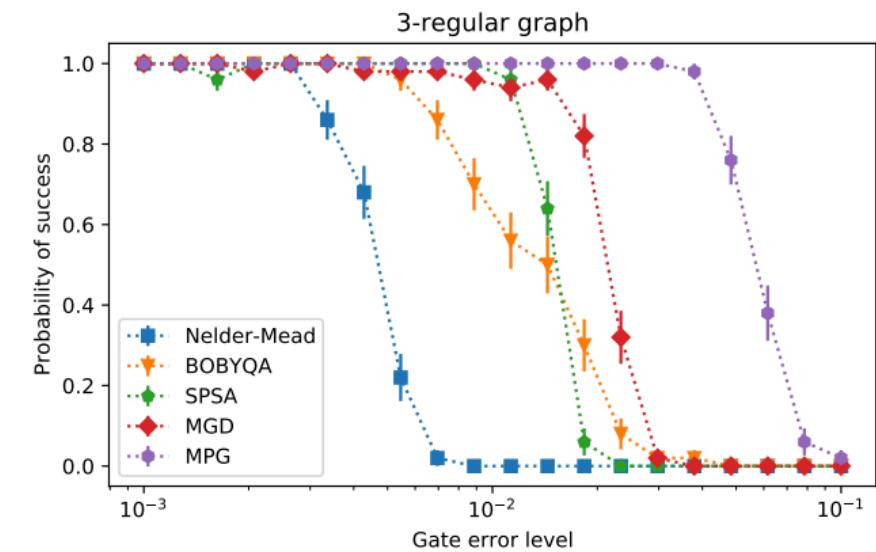
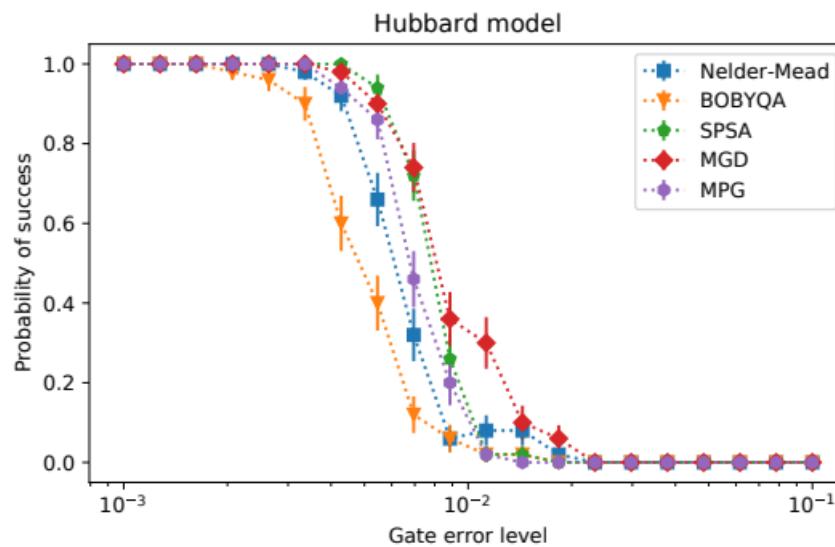
$$R_2(\boldsymbol{\theta}^{(T)}) \leq \tilde{O}\left(\text{poly}\left(\frac{d}{BK^2(1-p)^L}\right)\right).$$

Main results of the trainability of QNN towards ERM

JDT 京东科技

Implications of Theorem 1:

- $K \uparrow, B \uparrow, p \downarrow, d \downarrow$, and $L \downarrow$ yield better utility bounds $R_1(\theta^{(T)})$ and $R_2(\theta^{(T)})$;

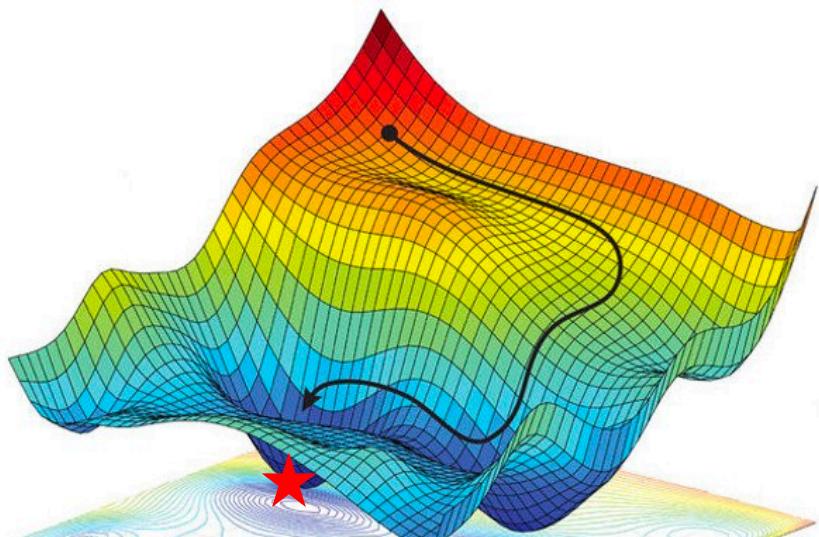


Main results of the trainability of QNN towards ERM

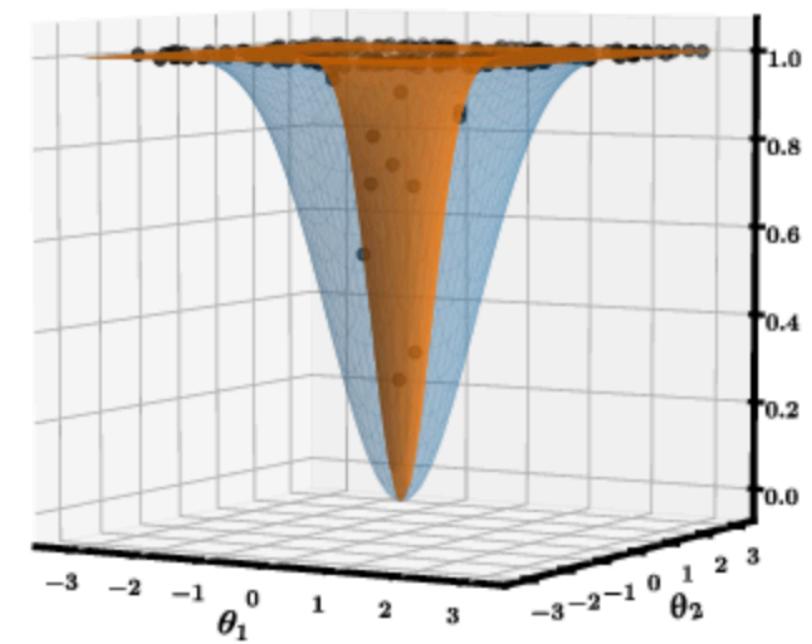
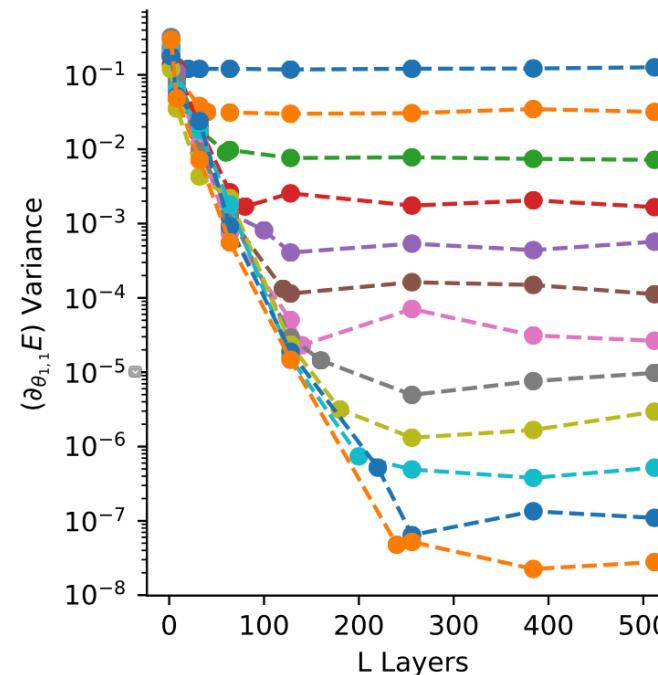
JDT 京东科技

Implications of Theorem 1:

- $K \uparrow, B \uparrow, p \downarrow, d \downarrow$, and $L \downarrow$ yield better utility bounds $R_1(\theta^{(T)})$ and $R_2(\theta^{(T)})$;
- The convergence towards the global optima as shown in $R_2(\theta^{(T)})$ provides an insight about how to employ regularization techniques to avoid barren plateaus.

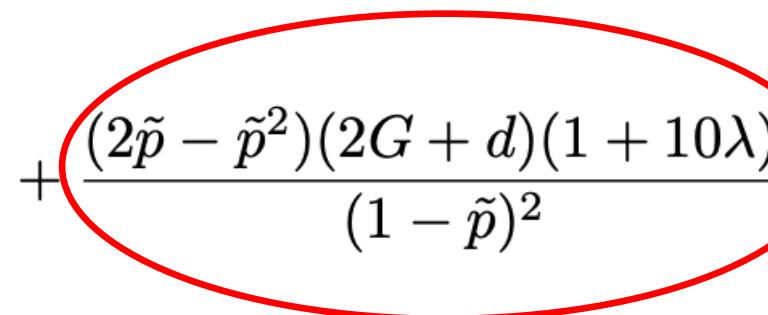


$R_2(\theta^{(T)})$ concerns the global minima,
the target result of QNN



Implications of Theorem 1:

- $K \uparrow, B \uparrow, p \downarrow, d \downarrow$, and $L \downarrow$ yield better utility bounds $R_1(\boldsymbol{\theta}^{(T)})$ and $R_2(\boldsymbol{\theta}^{(T)})$;
- The convergence towards the global optima as shown in $R_2(\boldsymbol{\theta}^{(T)})$ provides an insight about how to employ regularization techniques to avoid barren plateaus.
- The utility bound $R_1(\boldsymbol{\theta}^{(T)})$ indicates that the optimization of QNN can diverge for large d and p , no matter how large T or K would become;

$$R_1 \leq \frac{2S(1 + 90\lambda d)}{T(1 - \tilde{p})^2} + \frac{(2\tilde{p} - \tilde{p}^2)(2G + d)(1 + 10\lambda)^2}{(1 - \tilde{p})^2} + \frac{6dK + 8d}{(1 - \tilde{p})^2 BK^2}.$$


Sketch of proof

$$R_1 \leq \tilde{O} \left(\text{poly} \left(\frac{d}{T(1-p)^{L_Q}}, \frac{d}{BK(1-p)^{L_Q}}, \frac{d}{(1-p)^{L_Q}} \right) \right) \quad R_1(\boldsymbol{\theta}^{(T)}) := \mathbb{E} \left[\left\| \nabla \mathcal{L}(\boldsymbol{\theta}^{(T)}) \right\|^2 \right]$$

The objective function $\mathcal{L}(\boldsymbol{\theta})$ is S-smooth, i.e.,

$$\mathcal{L}(\boldsymbol{\theta}^{(t+1)}) - \mathcal{L}(\boldsymbol{\theta}^{(t)}) \leq \langle \nabla \mathcal{L}(\boldsymbol{\theta}^{(t)}), \boldsymbol{\theta}^{(t+1)} - \boldsymbol{\theta}^{(t)} \rangle + \frac{S}{2} \|\boldsymbol{\theta}^{(t+1)} - \boldsymbol{\theta}^{(t)}\|^2.$$

The optimization rule of **noisy** QNN at the t-th iteration follows

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta \nabla \bar{\mathcal{L}}(\boldsymbol{\theta}^{(t)}).$$

$$\mathcal{L}(\boldsymbol{\theta}^{(t+1)}) - \mathcal{L}(\boldsymbol{\theta}^{(t)})$$

$$\leq -\frac{1}{S} \langle \nabla \mathcal{L}(\boldsymbol{\theta}^{(t+1)}), \nabla \bar{\mathcal{L}}(\boldsymbol{\theta}^{(t)}) \rangle$$

$$+ \frac{1}{2S} \|\nabla \bar{\mathcal{L}}(\boldsymbol{\theta}^{(t)})\|^2$$

Sketch of proof

$$R_1 \leq \tilde{O} \left(\text{poly} \left(\frac{d}{T(1-p)^{L_Q}}, \frac{d}{BK(1-p)^{L_Q}}, \frac{d}{(1-p)^{L_Q}} \right) \right) \quad R_1(\boldsymbol{\theta}^{(T)}) := \mathbb{E} \left[\left\| \nabla \mathcal{L}(\boldsymbol{\theta}^{(T)}) \right\|^2 \right]$$

$$\begin{aligned} & \mathcal{L}(\boldsymbol{\theta}^{(t+1)}) - \mathcal{L}(\boldsymbol{\theta}^{(t)}) \\ & \leq -\frac{1}{S} \langle \nabla \mathcal{L}(\boldsymbol{\theta}^{(t+1)}), \nabla \bar{\mathcal{L}}(\boldsymbol{\theta}^{(t)}) \rangle \\ & \quad + \frac{1}{2S} \|\nabla \bar{\mathcal{L}}(\boldsymbol{\theta}^{(t)})\|^2 \end{aligned}$$

Theorem 3. It follows that

$$\nabla_j \bar{\mathcal{L}}_i(\boldsymbol{\theta}^{(t)}) = (1-\tilde{p})^2 \nabla_j \mathcal{L}_i(\boldsymbol{\theta}^{(t)}) + C_{j,1}^{(i,t)} + \xi_i^{(t,j)},$$

where $\tilde{p} = 1 - (1-p)^{L_Q}$, L_Q is the circuit depth, the constant $C_{j,1}^{(i,t)}$ only depends on Y_i , $\boldsymbol{\theta}^{(t)}$, and \tilde{p} , and $\xi_i^{(t,j)}$ follows the distribution \mathcal{P}_Q that is formed by Y_i , $\boldsymbol{\theta}^{(t)}$, the number of measurements K , and \tilde{p} with zero mean.

$$\mathbb{E}_{\xi_i^{(t)}, \xi_i^{(t,j)}} [\mathcal{L}(\boldsymbol{\theta}^{(t+1)}) - \mathcal{L}(\boldsymbol{\theta}^{(t)})] \leq -\frac{1}{2S} (1-\tilde{p})^2 \|\nabla \mathcal{L}(\boldsymbol{\theta}^{(t)})\|^2 + \frac{2G+d}{2S} (2-\tilde{p}) \tilde{p} (1+10\lambda)^2 + \frac{6dK+8d}{2SBK^2}.$$

Sketch of proof

$$R_1 \leq \tilde{O} \left(\text{poly} \left(\frac{d}{T(1-p)^{L_Q}}, \frac{d}{BK(1-p)^{L_Q}}, \frac{d}{(1-p)^{L_Q}} \right) \right) \quad R_1(\boldsymbol{\theta}^{(T)}) := \mathbb{E} \left[\left\| \nabla \mathcal{L}(\boldsymbol{\theta}^{(T)}) \right\|^2 \right]$$

$$\mathbb{E}_{\xi_i^{(t)}, \xi_i^{(t,j)}} [\mathcal{L}(\boldsymbol{\theta}^{(t+1)}) - \mathcal{L}(\boldsymbol{\theta}^{(t)})] \leq -\frac{1}{2S}(1-\tilde{p})^2 \left\| \nabla \mathcal{L}(\boldsymbol{\theta}^{(t)}) \right\|^2 + \frac{2G+d}{2S}(2-\tilde{p})\tilde{p}(1+10\lambda)^2 + \frac{6dK+8d}{2SBK^2}.$$

By induction, with summing over $t = 0, \dots, T-1$ and taking expectation, we obtain

$$\mathbb{E}_t \left[\left\| \nabla \mathcal{L}(\boldsymbol{\theta}^{(t)}) \right\|^2 \right] \leq \frac{2S(1+90\lambda d)}{T(1-\tilde{p})^2} + \frac{(2\tilde{p}-\tilde{p}^2)(2G+d)(1+10\lambda)^2}{(1-\tilde{p})^2} + \frac{6dK+8d}{(1-\tilde{p})^2 BK^2}.$$

Sketch of proof

JDT 京东科技

$$R_2 \leq (1 + 90\lambda d) \exp\left(-\frac{\mu(1 - \tilde{p})^2 T}{S}\right) + T \frac{(2\tilde{p} - \tilde{p}^2)(G + 2d)(1 + 10\lambda)^2 BK^2 + 6dK + 8d}{2SBK^2}$$

$$R_2(\boldsymbol{\theta}^{(T)}) := \mathbb{E}[\mathcal{L}(\boldsymbol{\theta}^{(T)})] - \mathcal{L}(\boldsymbol{\theta}^*),$$

PL condition: a function f satisfies PL condition if there exists $\mu > 0$ and for every possible $\boldsymbol{\theta} \in \mathcal{C}$, $\|\nabla f(\boldsymbol{\theta})\|^2 \geq 2\mu(f(\boldsymbol{\theta}) - f^*)$, where $f^* = \min_{\boldsymbol{\theta} \in \mathcal{C}} f(\boldsymbol{\theta})$.

Lemma Assume $\lambda \in \left(0, \frac{1}{3\pi}\right) \cup \left(\frac{1}{\pi}, \infty\right)$ the loss \mathcal{L} used in QNN satisfies PL condition with $\mu = \frac{(-1 + \lambda\pi)^2}{1 + \lambda d(3\pi)^2}$.

Sketch of proof

$$R_2 \leq (1 + 90\lambda d) \exp\left(-\frac{\mu(1 - \tilde{p})^2 T}{S}\right) + T \frac{(2\tilde{p} - \tilde{p}^2)(G + 2d)(1 + 10\lambda)^2 BK^2 + 6dK + 8d}{2SBK^2}$$

$$R_2(\boldsymbol{\theta}^{(T)}) := \mathbb{E}[\mathcal{L}(\boldsymbol{\theta}^{(T)})] - \mathcal{L}(\boldsymbol{\theta}^*),$$

Lemma Assume $\lambda \in \left(0, \frac{1}{3\pi}\right) \cup \left(\frac{1}{\pi}, \infty\right)$ the loss \mathcal{L} used in QNN satisfies PL condition with

$$\mu = \frac{(-1 + \lambda\pi)^2}{1 + \lambda d(3\pi)^2}.$$

$$+ \quad \|\nabla f(\boldsymbol{\theta})\|^2 \geq 2\mu(f(\boldsymbol{\theta}) - f^*)$$

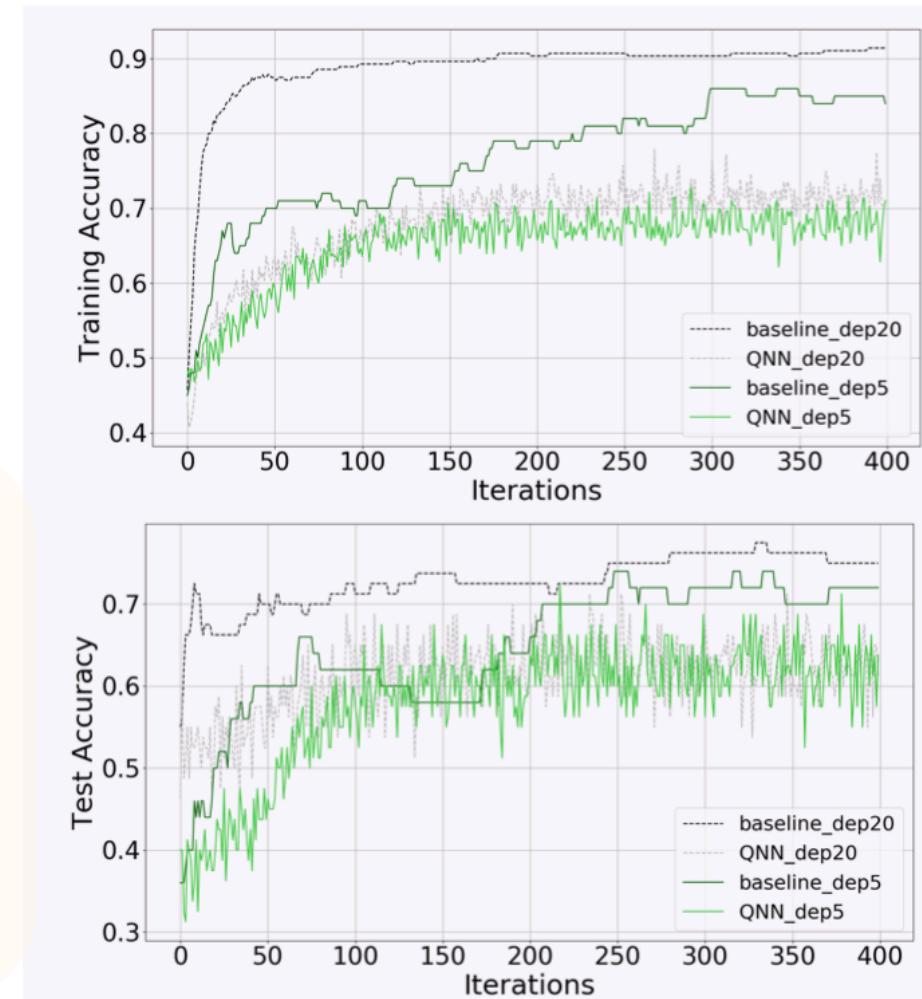
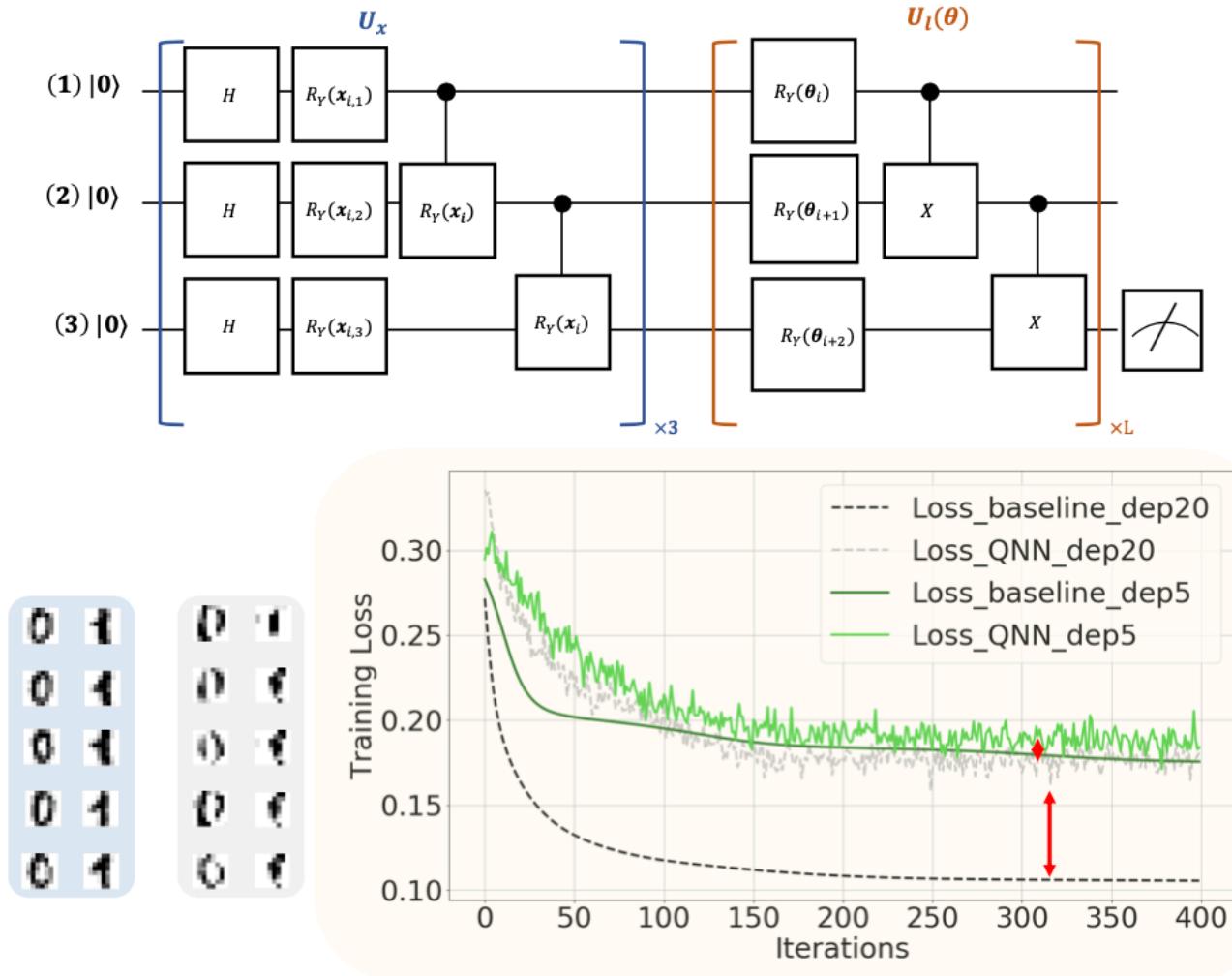
$$\mathbb{E}_{\xi_i^{(t)}, \xi_i^{(t,j)}} [\mathcal{L}(\boldsymbol{\theta}^{(t+1)}) - \mathcal{L}(\boldsymbol{\theta}^{(t)})] \leq -\frac{1}{2S}(1 - \tilde{p})^2 \|\nabla \mathcal{L}(\boldsymbol{\theta}^{(t)})\|^2 + \frac{2G + d}{2S}(2 - \tilde{p})\tilde{p}(1 + 10\lambda)^2 + \frac{6dK + 8d}{2SBK^2}.$$

By induction, with summing over t, and taking expectation

$$\mathbb{E}_{\boldsymbol{\varsigma}^{(t)}} [\mathcal{L}(\boldsymbol{\theta}^{(T)})] - \mathcal{L}^* = \frac{(-1 + \lambda\pi)^2}{1 + \lambda d(3\pi)^2}.$$

Lemma Assume $\lambda \in \left(0, \frac{1}{3\pi}\right) \cup \left(\frac{1}{\pi}, \infty\right)$ the loss \mathcal{L} used in QNN satisfies PL condition with

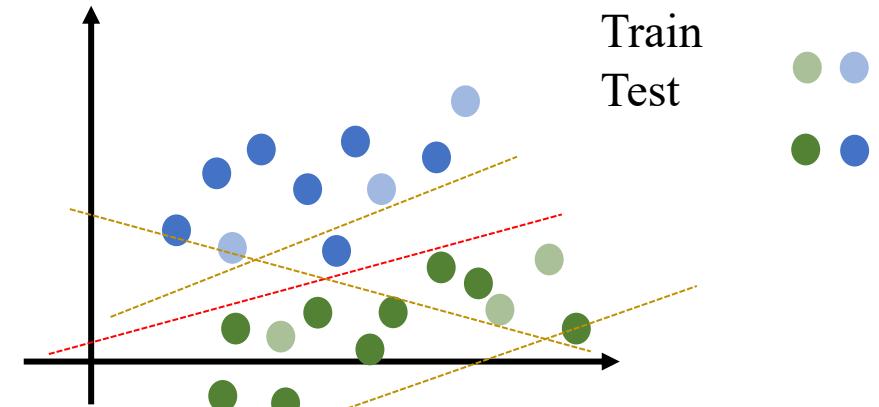
Simulation results



Problem setup of the generalization of QNN

- **Generalization**: for a given learning problem, QNN requires a poly sample complexity to find a good hypothesis, while classical algorithms require an exp number of samples.

(Quantum) Learning theory



Quantum statistical query (QSQ) oracle: A QSQ oracle which takes a tolerance parameter τ and an observable $\mathbb{M} \in \mathbb{C}^{2^{N+1} \times 2^{N+1}}$ and returns a number α satisfies

$$|\alpha - \langle \psi_{c^*} | \mathbb{M} | \psi_{c^*} \rangle| \leq \tau.$$

$C \subseteq \{c: \{0,1\}^N \rightarrow \{0,1\}\}$: a concept class;

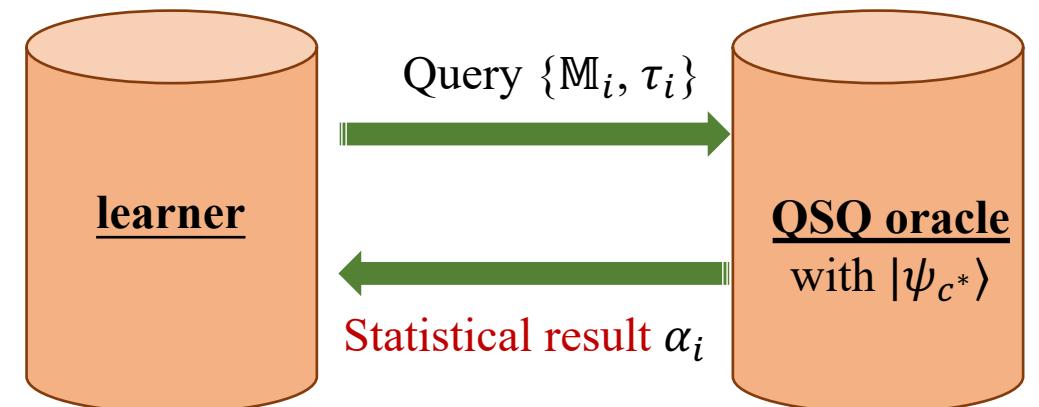
$D: \{0,1\}^N \rightarrow [0, 1]$: an unknown distribution;

$|\psi_{c^*}\rangle = \sum_{x \in \{0,1\}^N} \sqrt{D(x)} |x, c^*(x)\rangle$: a quantum example.

Central tools to explore the generalization of QNN

Quantum statistical query (QSQ) learning algorithm: The QSQ learning algorithm adaptively feeds a sequence of $\{\mathbb{M}_i, \tau_i\}$ into a QSQ oracle, and exploits the responses of $\{\alpha_i\}$ to output a hypothesis $h: \{0,1\}^N \rightarrow \{0,1\}$.

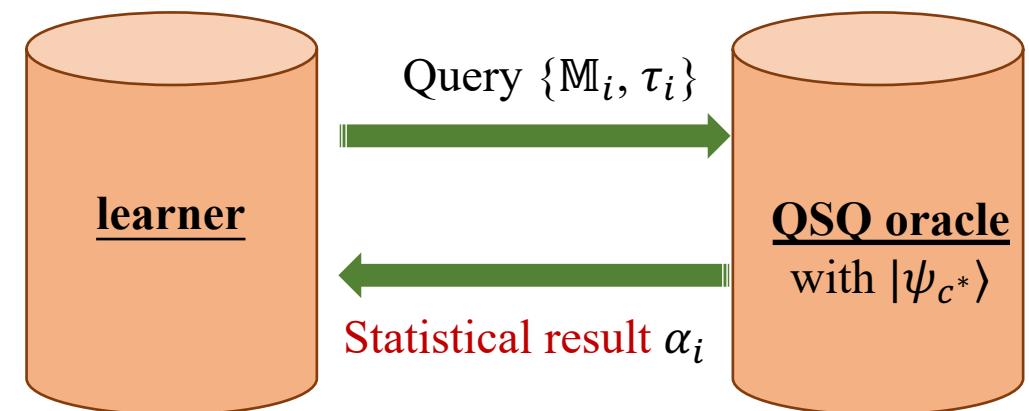
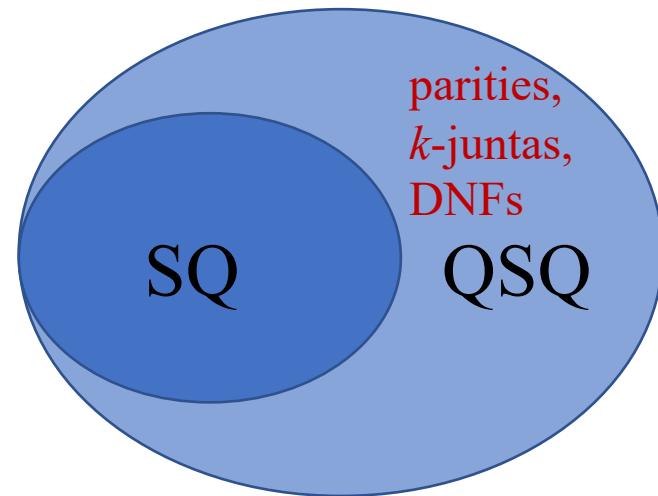
The goal of the learner is to achieve $\Pr_{x \sim D}(h(x) = c^*(x)) \leq \epsilon$ for all possible D and c^* .



Central tools to explore the generalization of QNN

JDT 京东科技

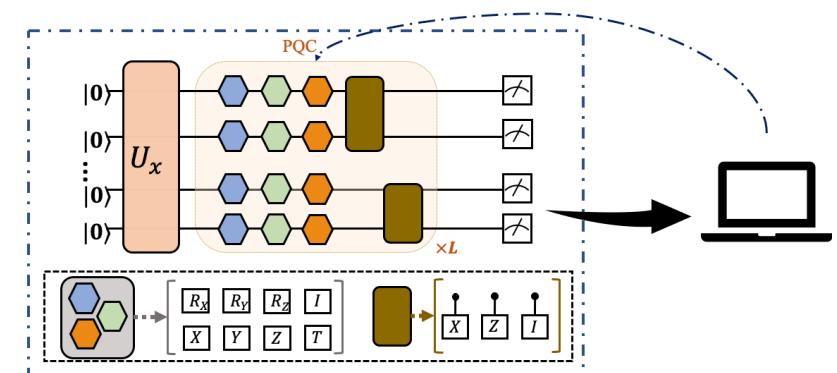
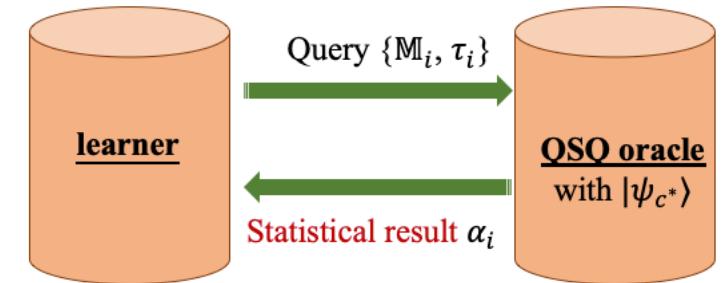
Lemma 1. Let C be the concept class of parities, k -juntas, or $\text{poly}(N)$ -sized DNFs (Disjunctive Normal Forms), then there exists a $\text{poly}(N)$ -query QSQ algorithm with tolerance $\tau = \tilde{O}(\epsilon)$ that efficiently learns C under the uniform distribution.



The generalization of QNN

Motivation: the QSQ oracle yields a similar behavior of the variational quantum circuit used in QNN:

- 1) measurement statistics
- 2) noisy setting.



Main results of the generalization of QNN

JDT 京东科技

Theorem 2. A QSQ learning algorithm, where the distribution over the quantum example $|\psi_{c^*}\rangle$ is fixed to be uniform and the observable M can be implemented by at most $\text{poly}(n)$ single and two-qubit gates, can be **efficiently simulated by noisy QNN** using **polynomial** samples.

Main results of the generalization of QNN

JDT 京东科技

Theorem 2. A QSQ learning algorithm, where the distribution over the quantum example $|\psi_{c^*}\rangle$ is fixed to be uniform and the observable \mathbb{M} can be implemented by at most $\text{poly}(n)$ single and two-qubit gates, can be efficiently simulated by noisy QNN using **polynomial** samples.

Implications:

- The noisy quantum circuit used in QNN has potential to tackle practical learning tasks, e.g., **support vector machines** and **linear and convex optimization**, with quantum advantages.

QSQ oracle: A QSQ oracle which takes a tolerance parameter τ and an observable $\mathbb{M} \in \mathbb{C}^{2^{N+1} \times 2^{N+1}}$ and returns a number α satisfies

$$|\alpha - \langle \psi_{c^*} | \mathbb{M} | \psi_{c^*} \rangle| \leq \tau.$$


 ν

Let the encoding circuit U_x prepare the state $|\psi_{c^*}\rangle$ and the quantum measurement constructed from \mathbb{M} . Under the **depolarization noise**, the expectation value of quantum measurements of the noisy QNN yields

$$\tilde{\nu} = (1 - \tilde{p})\nu + \frac{\tilde{p} \operatorname{Tr}(\mathbb{M})}{2^{N+1}}$$

Let the encoding circuit U_x prepare the state $|\psi_{c^*}\rangle$ and the quantum measurement constructed from \mathbb{M} . Under the **depolarization noise**, the expectation value of quantum measurements of the noisy QNN yields

$$\tilde{\nu} = (1 - \tilde{p})\nu + \frac{\tilde{p} \text{Tr}(\mathbb{M})}{2^{N+1}}$$

+

By the Chernoff-Hoeffding bound for real-valued variables, the relation between the **sample mean** $\frac{1}{K} \sum_{i=1}^K V_k$ and the result $\tilde{\nu}$ follows

$$\Pr \left(\left| \frac{1}{K} \sum_{i=1}^K V_k - \tilde{\nu} \right| \geq \frac{\delta}{2} \right) \leq 2 \exp(-\delta^2 K / 2).$$



$$\left| \frac{1}{K} \sum_{k=1}^K V_k - \nu \right| \leq \tilde{p} \left(\nu + \frac{1}{2^{N+1}} \right) + \frac{\delta}{2} \leq \frac{5}{4} \tilde{p} + \frac{\delta}{2} \leq \tau$$



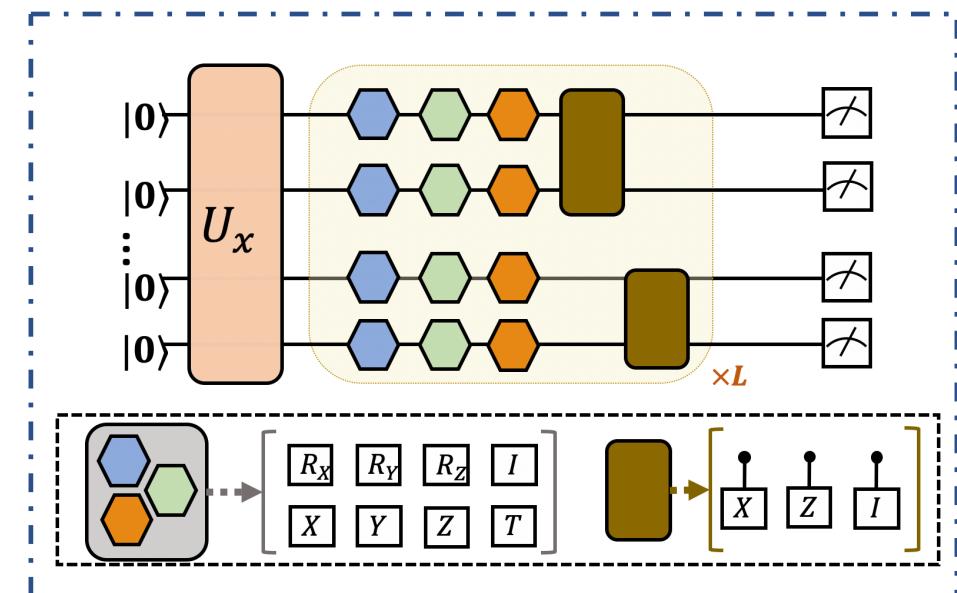
With the successful probability $1 - 2b$, noisy QNN can efficiently simulate QSQ oracle using $K = 2 \ln \left(\frac{2}{b} \right) / \delta^2$ examples.

- ✓ More measurements, lower noise, and shallower circuit depth contribute to a better performance of QNN.
 - ✓ Theoretical guidance to devise more advanced QNN based learning models that are robust to inevitable gate noise and insensitive to the barren plateau phenomenon.
 - ✓ QNN can efficiently learn parity, juntas, and DNF with quantum advantages even with gate noise.
-
- Exploit more advanced quantum models developed in quantum learning theory to explore the potential advantages that can be achieved by QNN.

Part II. QAS: An efficient scheme to enhance the trainability of QNN and suppress its error

The **caveats** of NISQ processors are:

- limited number of qubits;
- shallow circuit depth;
- connectivity restriction;
- unavoidable system noise.



Problem setup

Quantum neural networks (QNN) aims to find the optimal $\boldsymbol{\theta}^* \in \mathbb{R}^d$ that minimizes an objective function \mathcal{L} , i.e.,

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta} \in \mathcal{C}} \mathcal{L}(\boldsymbol{\theta}, \mathbf{z}) := \frac{1}{n} \sum_{j=1}^n l(y_j, \hat{y}_j) + r(\boldsymbol{\theta}).$$

$\mathbf{z} = \{\mathbf{z}_j\}_{j=1}^N \in \mathcal{Z}$: Input dataset with N samples;

$\mathbf{z}_j = (\mathbf{x}_j, y_j)$: The j-th sample, the feature vector as $\mathbf{x}_j \in \mathbb{R}^{D_c}$, the label as $y_j \in \mathbb{R}^1$,

$U(\boldsymbol{\theta})$: The trainable quantum circuit as.

$\hat{y}_j = \text{Tr}((\Pi |\psi(\boldsymbol{\theta})\rangle\langle\psi(\boldsymbol{\theta})|))$: The predicted label of the j -th sample;

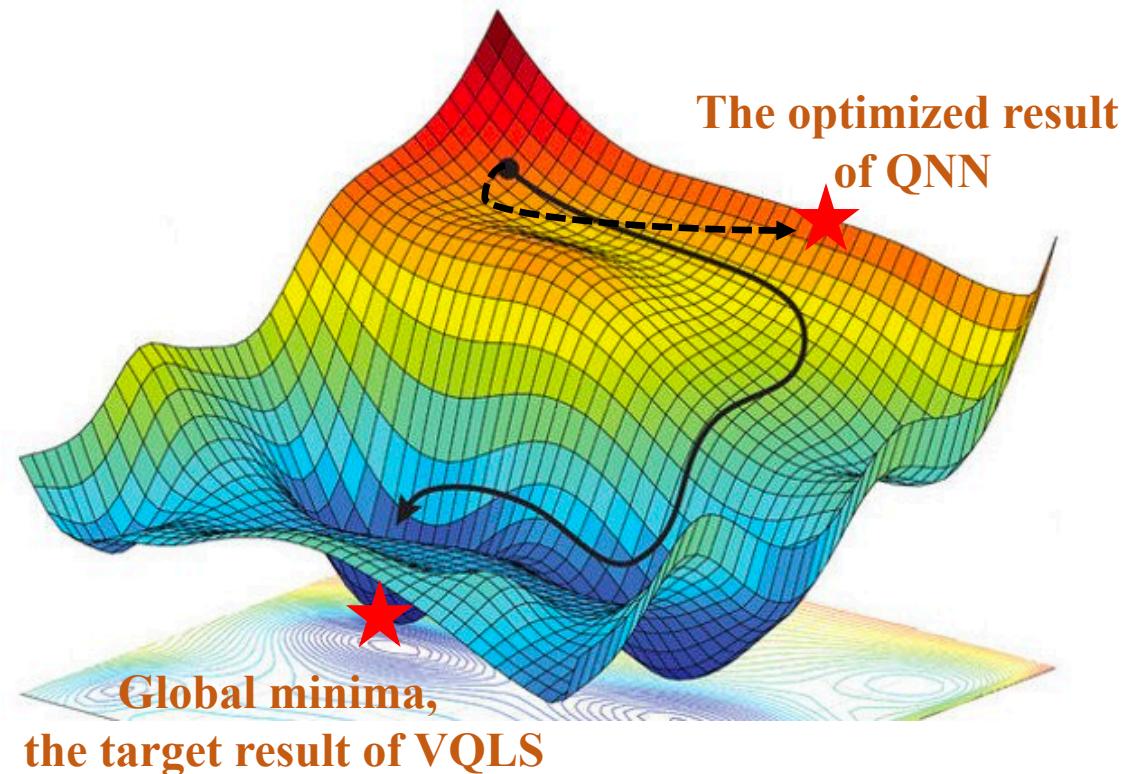
$|\psi(\boldsymbol{\theta})\rangle$: quantum states prepared by $U(\boldsymbol{\theta})$ and \mathbf{z}_j ;

$l(\cdot, \cdot)$: the loss function measures the difference between y_j and \hat{y}_j ;

$r(\boldsymbol{\theta})$: the regularizer.

Two key issues that destroy the optimization of QNN are:

- barren plateaus,
- accumulated system noise.



- To avoid barren plateaus:
 - ❖ modify cost functions (*Cerezo, M. et al. “Cost-Function-Dependent Barren Plateaus in Shallow Quantum Neural Networks.”*),
 - ❖ adopt some strategies to initialize parameters (*Grant, Edward, et al. "An initialization strategy for addressing barren plateaus in parametrized quantum circuits."*).
- To suppress quantum system noise:
 - ❖ quantum error mitigation techniques (*Endo, Suguru, et al. "Hybrid quantum-classical algorithms and quantum error mitigation."*).

Q: Can we unify the elimination of barren plateaus and the mitigation of the quantum errors in **a single problem?**

Can we unify the elimination of barren plateaus and the mitigation
of the quantum errors in **a single problem?**

Yes

Decreased circuit depth

The expressive power of $U(\theta)$ is decreased;
The barren plateaus' phenomenon is negligible;
The accumulated noise is negligible;
The optimized result $\theta^{(T)}$ is far away to θ^*

Increased circuit depth

The expressive power of $U(\theta)$ is improved;
The barren plateaus' phenomenon is serious;
The accumulated noise becomes larger;
The optimized result $\theta^{(T)}$ is far way from θ^*

Decreased circuit depth

The expressive power of $U(\theta)$ is decreased;
The barren plateaus' phenomenon is negligible;
The accumulated noise is negligible;
The optimized result $\theta^{(T)}$ is far away to θ^*



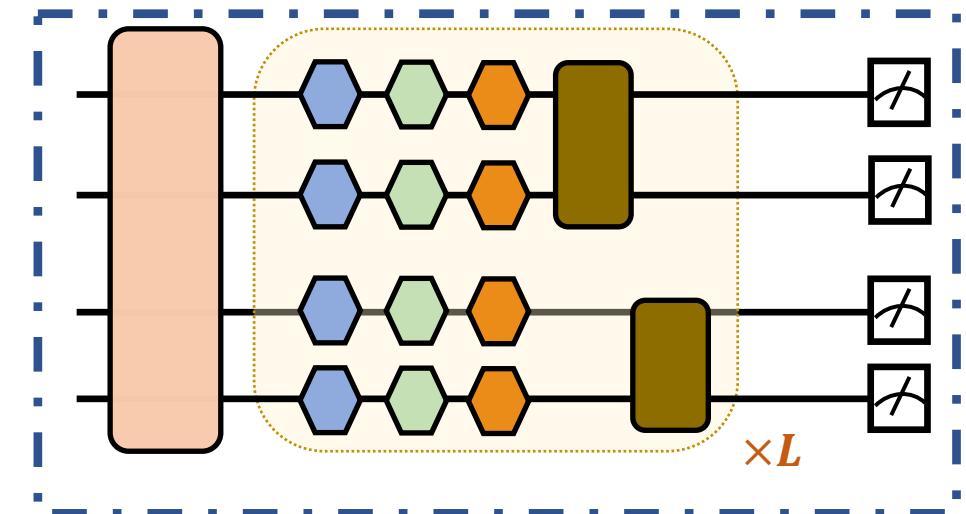
Target point

Increased circuit depth

The expressive power of $U(\theta)$ is sufficient;
The accumulated noise is tolerable;
The optimized result $\theta^{(T)}$ is close to θ^*

The elimination of barren plateaus and the suppression of quantum noise amount to seeking an appropriate circuit architecture of $U(\theta)$ satisfying the following requirements:

- a good expressive power to cover the target result θ^* ;
- using few number of quantum gates and shallow circuit depth;
- adapting to different connectivity constraints of various quantum processors.



Considering the circuit architecture information and the related noise, the objective of QNN is rewritten as

$$(\boldsymbol{\theta}^*, \mathbf{a}^*) = \arg \min_{\boldsymbol{\theta} \in \mathcal{C}, \mathbf{a} \in \mathcal{S}} \mathcal{L}(\boldsymbol{\theta}, \mathbf{z}, \mathbf{a}, \mathcal{E}_{\mathbf{a}})$$

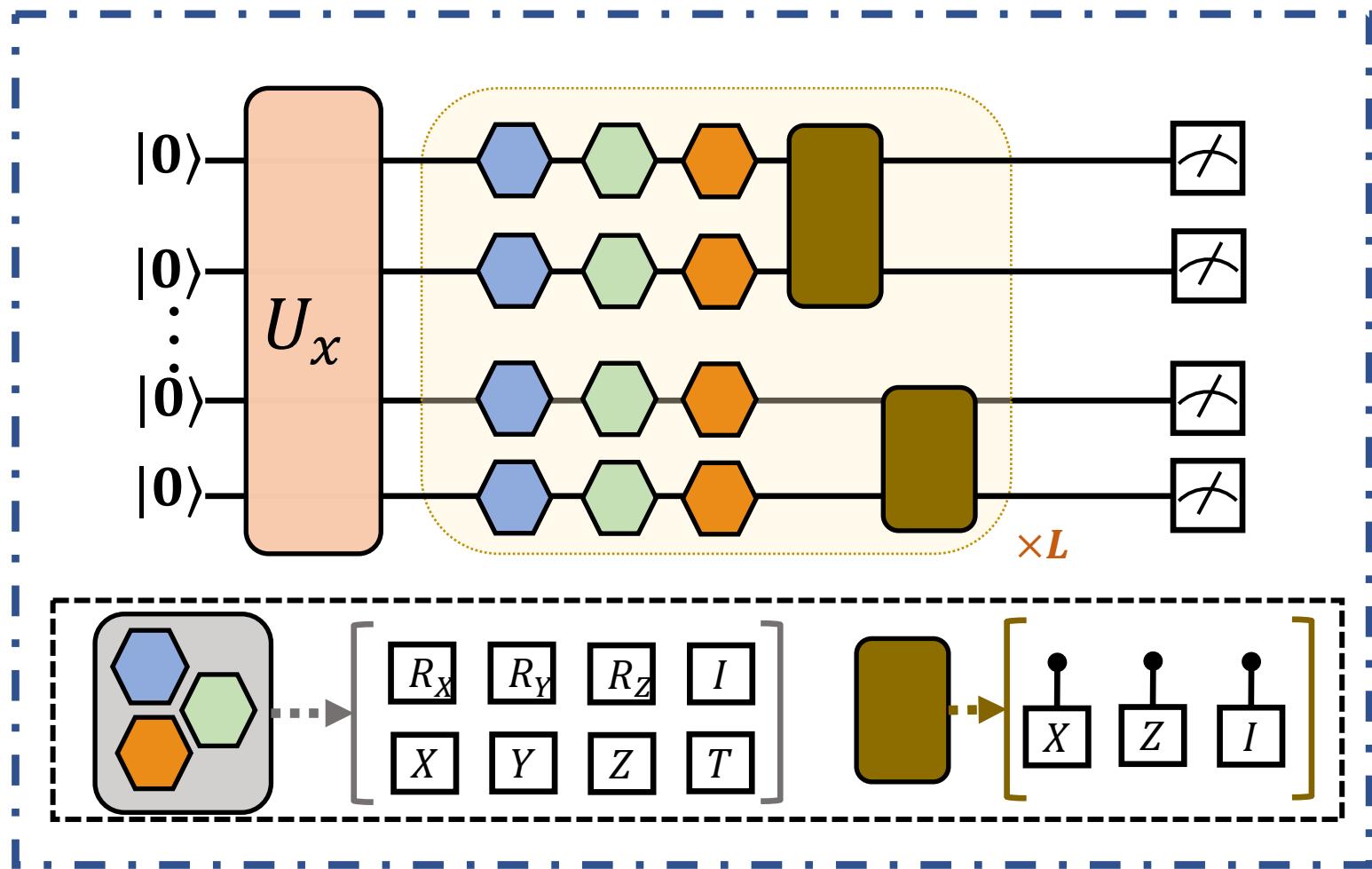
\mathcal{S} : the set that contains all possible circuit to build $U(\boldsymbol{\theta})$;

$\mathcal{E}_{\mathbf{a}}$: the quantum channel that simulates the quantum system noise induced by $\mathbf{a} \in \mathcal{S}$

Remarks

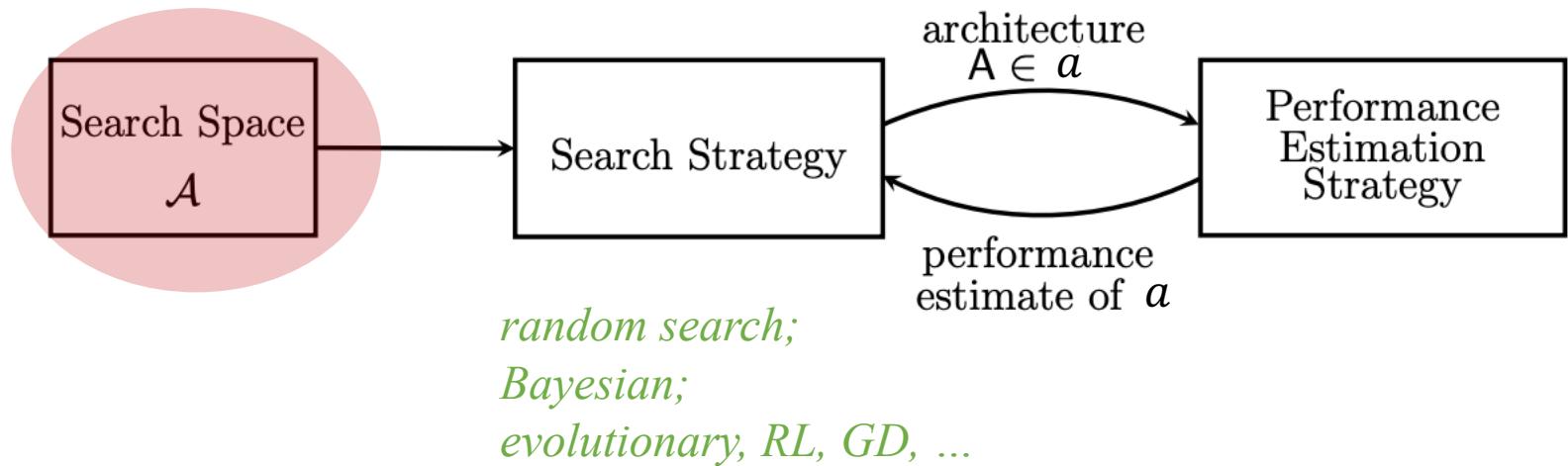
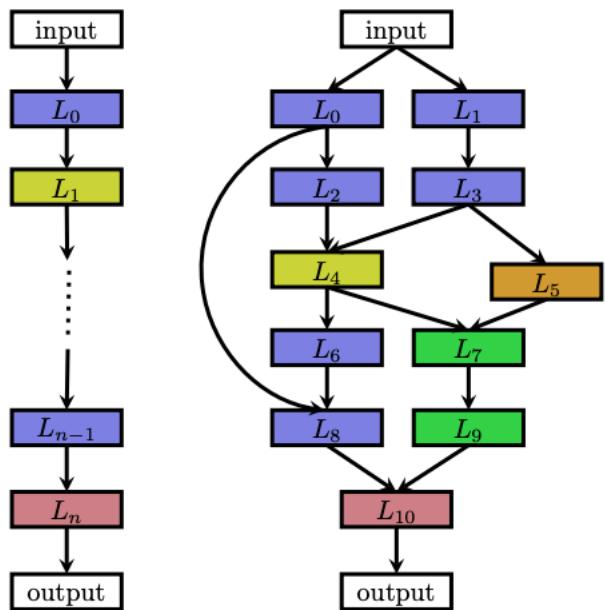
The brute-force searching is infeasible, since the total number of circuit architectures $|\mathcal{S}|$ **exponentially scales** with the qubits count N and the circuit depth L .

$$O((8^{3N} \times 3^N)^L)$$



AutoML: Neural network architecture search (NAS)

Different architecture spaces to search

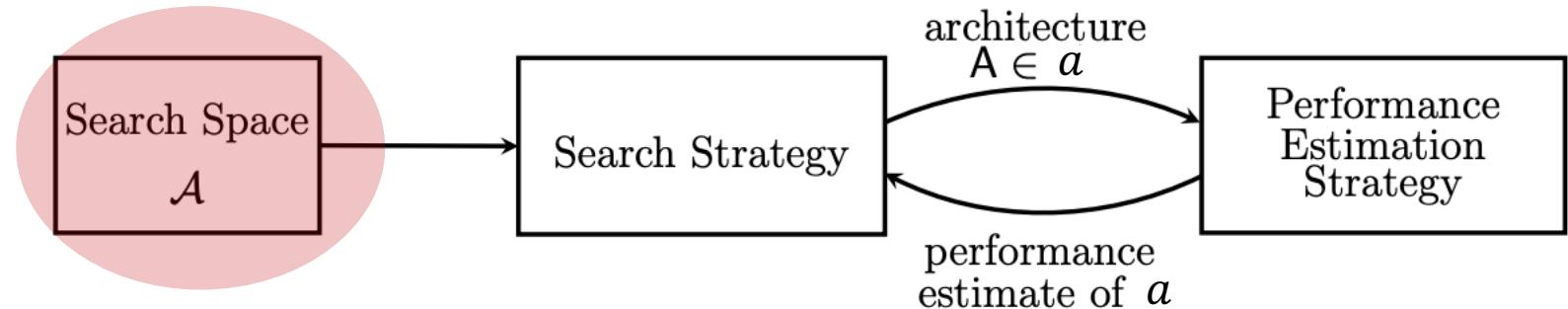


$$w_a = \underset{w}{\operatorname{argmin}} \mathcal{L}_{\text{train}} (\mathcal{N}(a, w))$$

$$a^* = \underset{a \in \mathcal{A}}{\operatorname{argmaxACC}_{\text{val}}} (\mathcal{N}(a, w_a))$$

Classical analog

Exponentially
large



$$w_a = \underset{w}{\operatorname{argmin}} \mathcal{L}_{\text{train}}(\mathcal{N}(a, w))$$

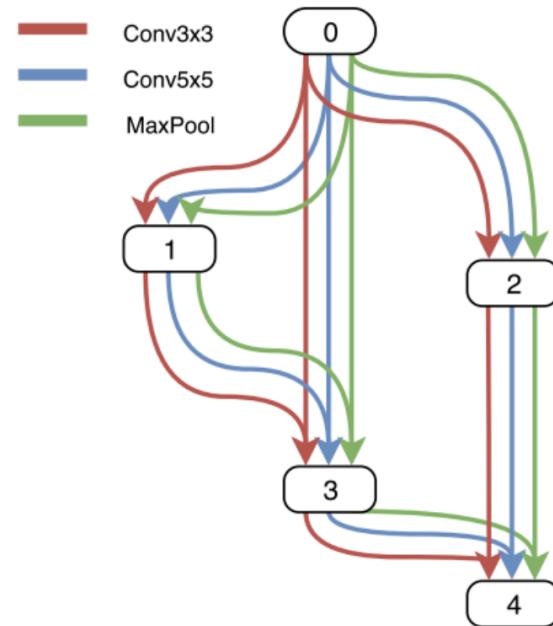
$$a^* = \underset{a \in \mathcal{A}}{\operatorname{argmaxACC}_{\text{val}}}(\mathcal{N}(a, w_a))$$

Optimization of NAS

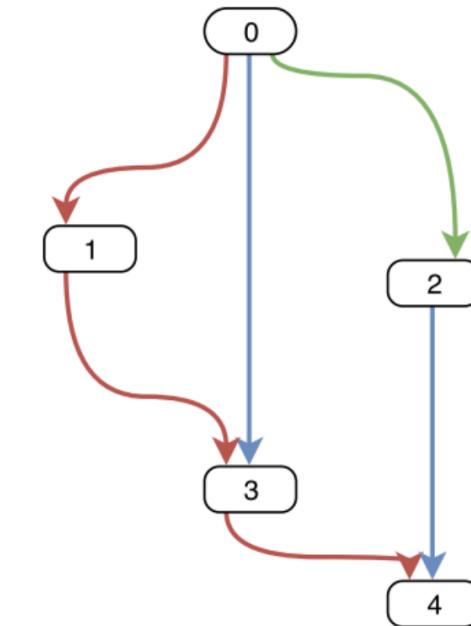
$$W_{\mathcal{A}} = \underset{W}{\operatorname{argmin}} \mathcal{L}_{\text{train}}(\mathcal{N}(\mathcal{A}, W))$$

$$a^* = \underset{a \in \mathcal{A}}{\operatorname{argmaxACC}_{\text{val}}}(\mathcal{N}(a, W_{\mathcal{A}}(a)))$$

A strategy to avoid exponential overhead: Weight Sharing



Supergraph / supernet \mathcal{A}



subgraph / subnet a

Quantum circuit architecture search (QAS)

JDT 京东科技

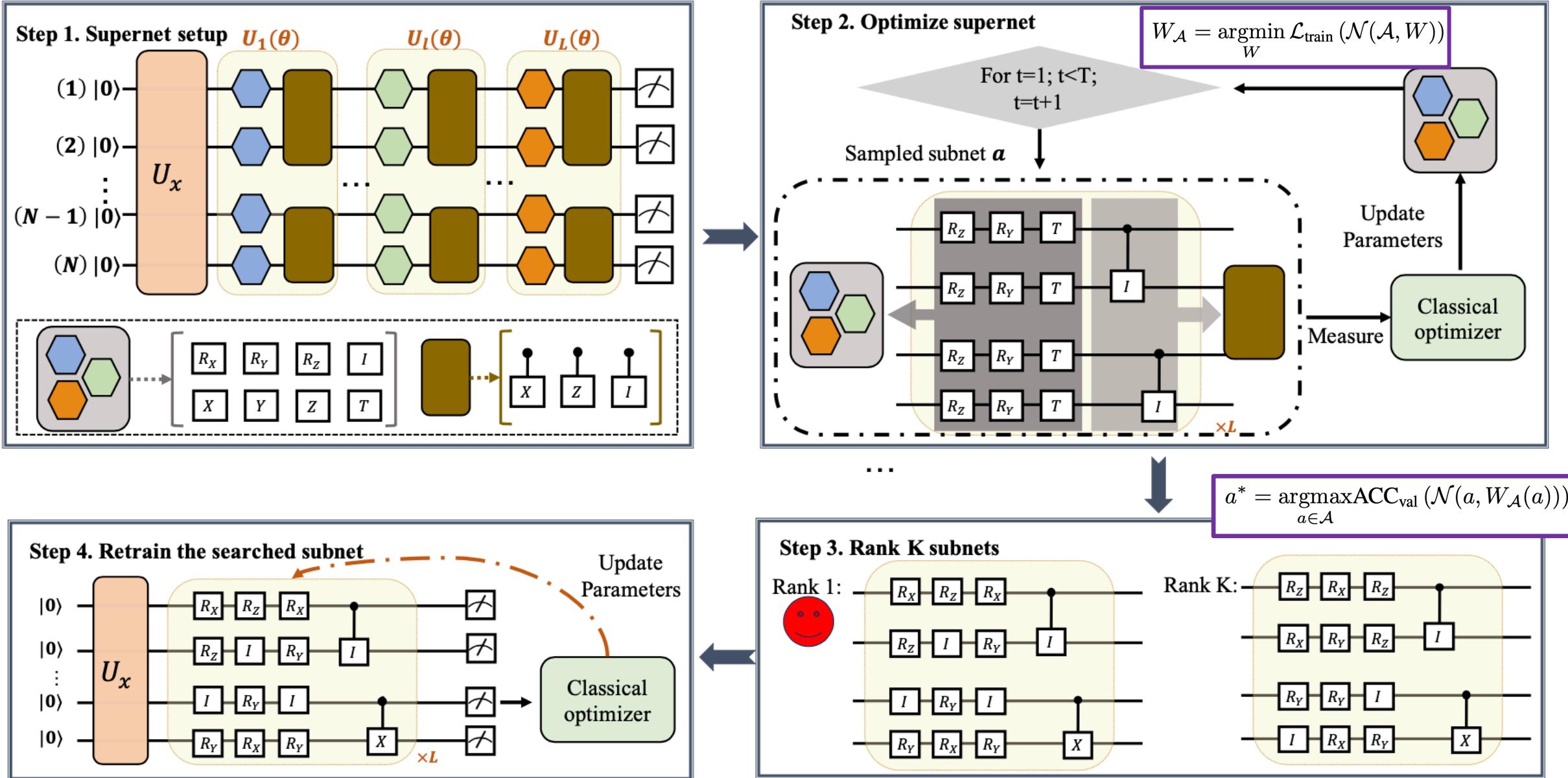
Here we propose the quantum circuit architecture search scheme (QAS) to tackle

$$(\boldsymbol{\theta}^*, \mathbf{a}^*) = \arg \min_{\boldsymbol{\theta} \in \mathcal{C}, \mathbf{a} \in \mathcal{S}} \mathcal{L}(\boldsymbol{\theta}, \mathbf{z}, \mathbf{a}, \mathcal{E}_a).$$

- The output of QAS **avoids barren plateaus** and **suppresses errors**;
- The runtime cost of QAS is **almost same with** conventional QNN-based algorithms.

The paradigm of QAS

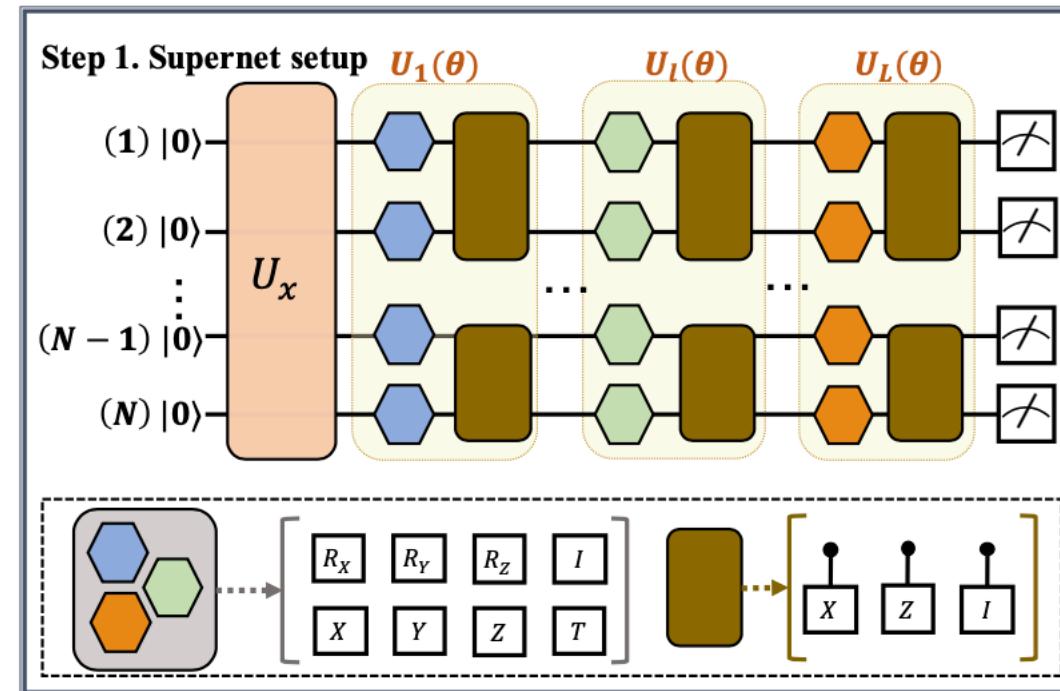
JDT 京东科技



Step 1 of QAS: supernet setup + weight sharing strategy

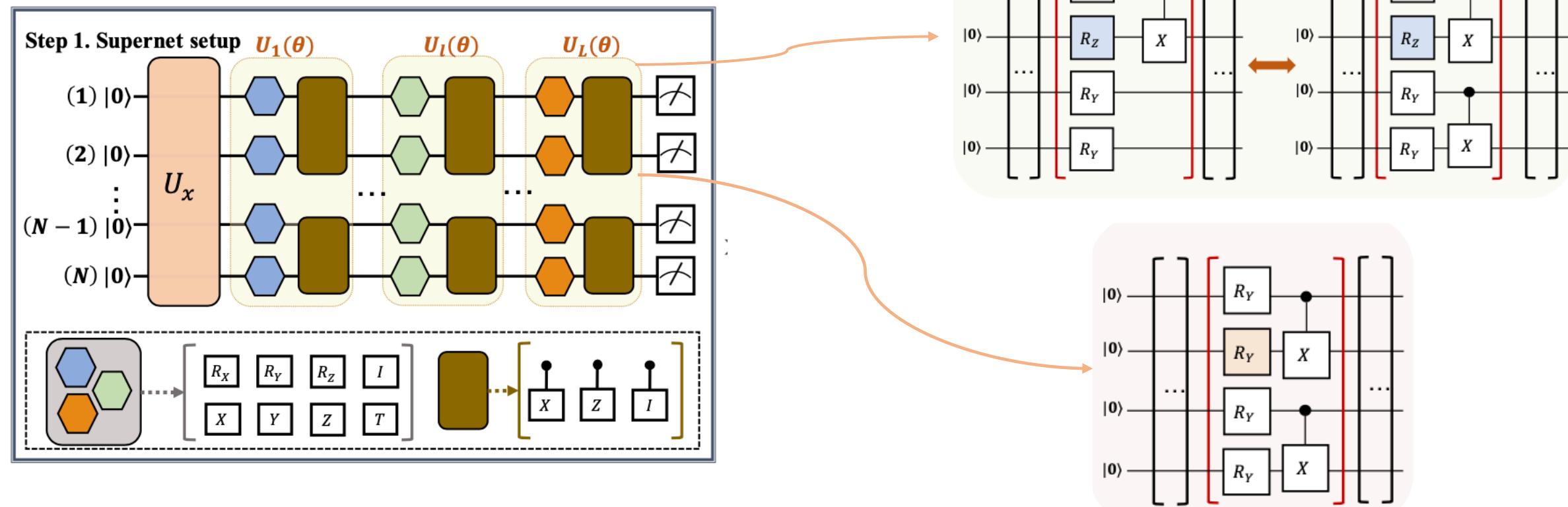
JDT 京东科技

The aim of the supernet in QAS is **manipulating all possible variational quantum circuit architectures of $U(\theta)$.**



Step 1 of QAS: supernet setup + weight sharing strategy

The weight sharing strategy: the weights in the supernet are shared across different architectures under the block level.

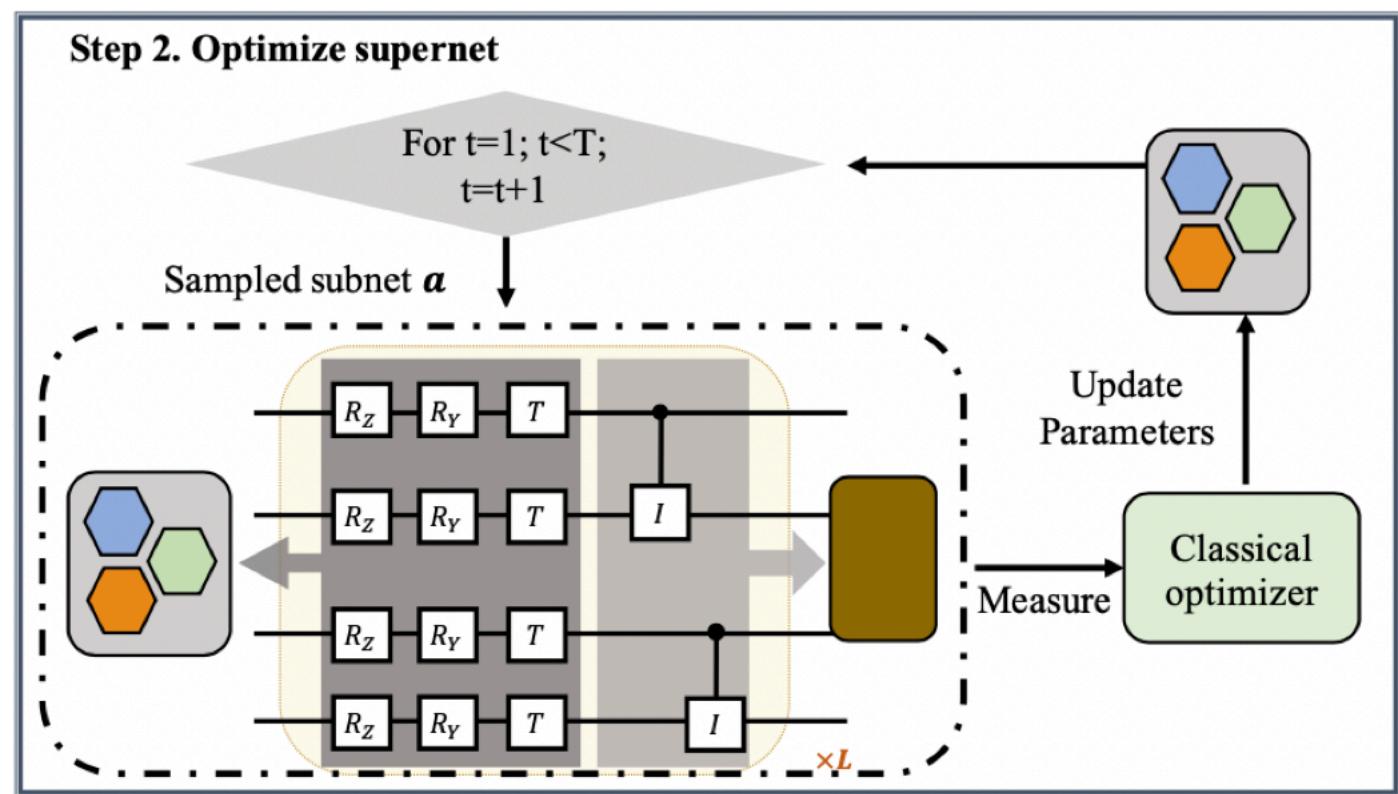
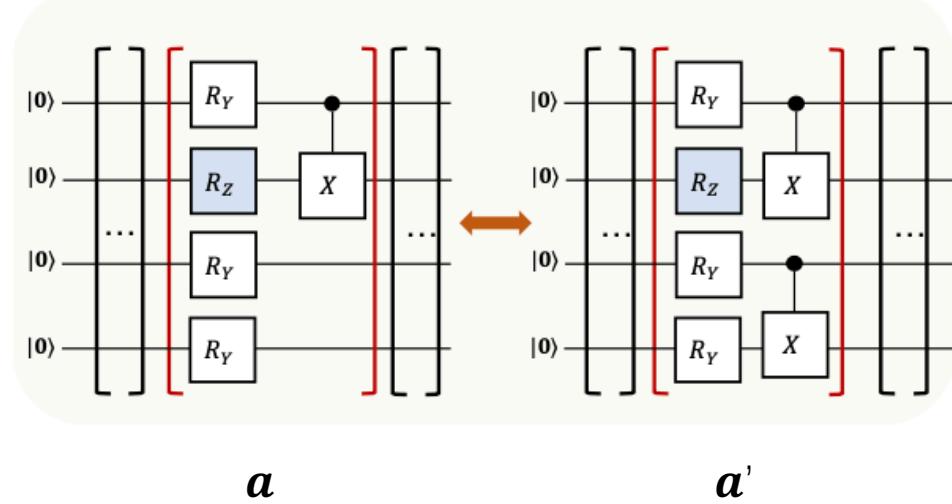


Step 2 of QAS: optimize supernet

At the t -th iteration, QAS uniformly samples a subnet $\mathbf{a}^{(t)} \in \mathcal{S}$ to minimize $\mathcal{L}(\boldsymbol{\theta}, \mathbf{z}, \mathbf{a}, \mathcal{E}_\mathbf{a})$, i.e., the parameters of the supernet relating to $\mathbf{a}^{(t)}$ is updated to

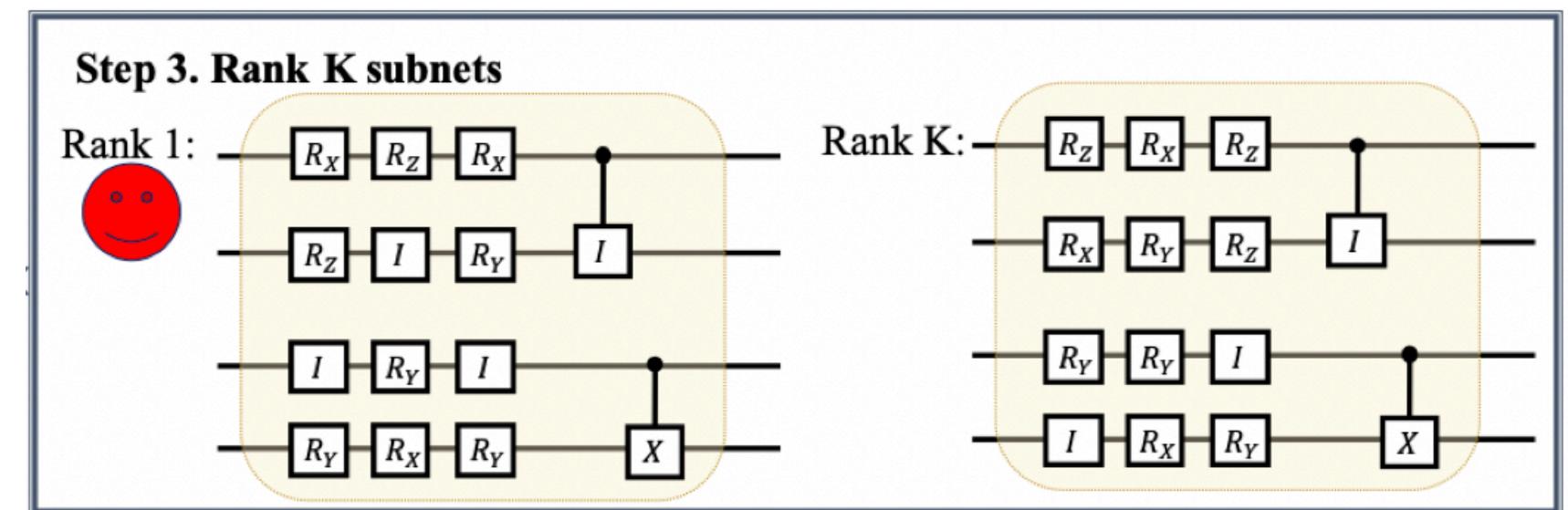
$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} + \eta \frac{\partial \mathcal{L}(\boldsymbol{\theta}^{(t)}, \mathbf{z}, \mathbf{a}, \mathcal{E}_\mathbf{a})}{\partial \boldsymbol{\theta}^{(t)}},$$

where η is the learning rate.



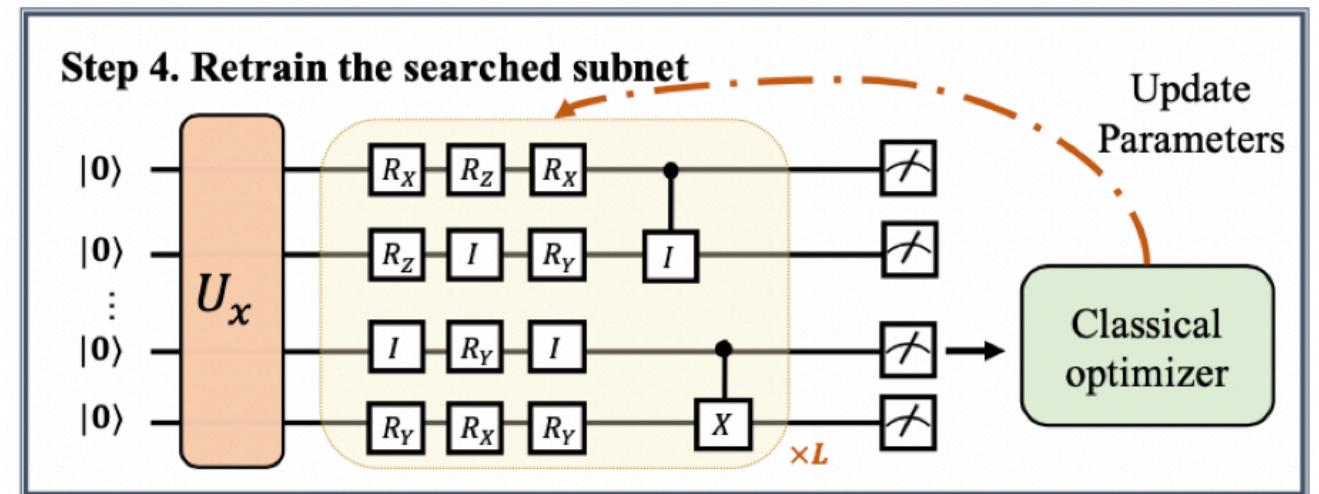
Step 3 of QAS: Rank K subnet

QAS uniformly samples K subnets, then ranks their performances, and assigns the subnet with the best performance as the output to approximate \mathbf{a}^* .



Step 4 of QAS: Rank K subnet

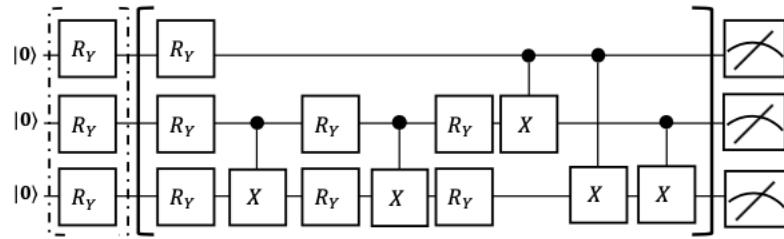
QAS retrain the searched subnet to approximate the optimal parameters.



The fierce competition phenomenon of QAS

Fierce competition in QAS: the performance of a specified subnet may be largely varied by independently training and QAS.

1st subnet a_1

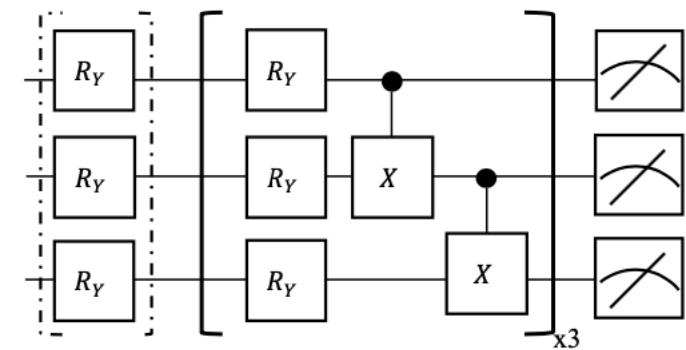


Independent training

99% acc | 94% acc

QAS with the weight sharing
91% acc | 93% acc

2cd subnet a_2

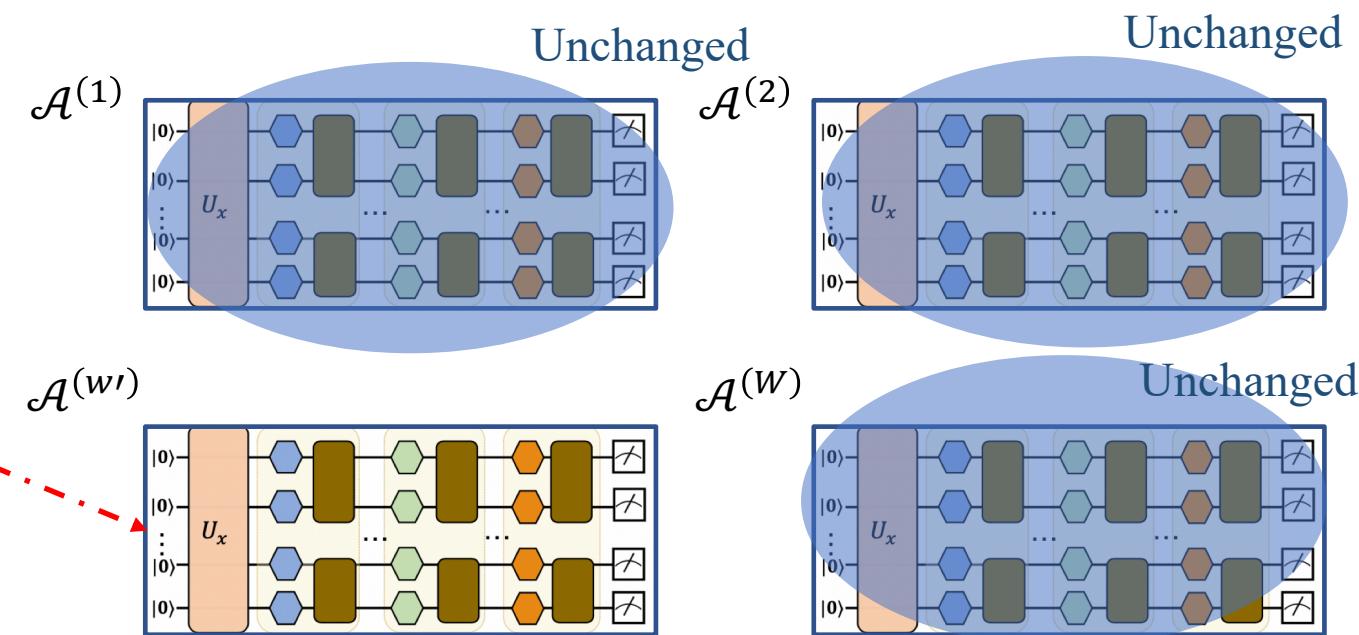
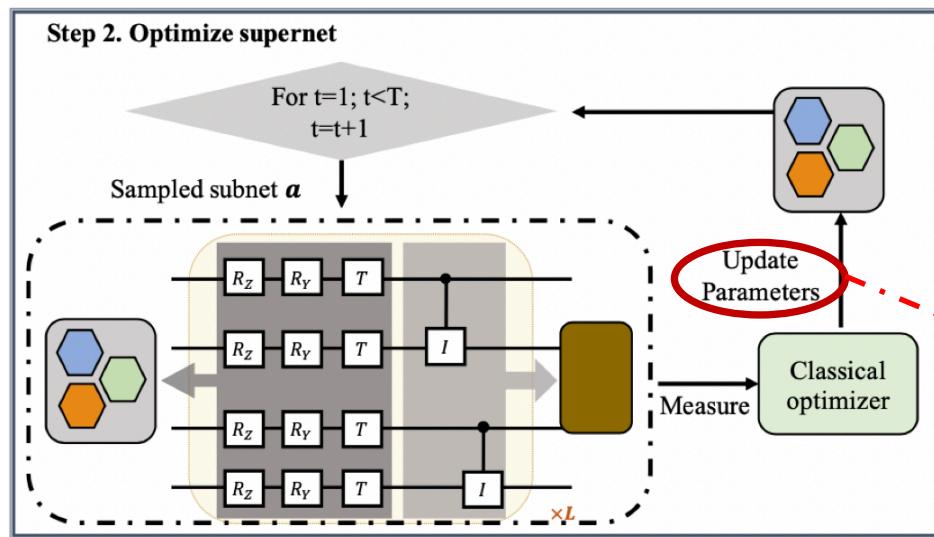


Solution to the fierce competition phenomenon of QAS

JDT 京东科技

Solution: Employ multiple supernets $\{\mathcal{A}^{(w)}\}_{w=1}^W$ and optimize them independently. In the stages II-III, the sampled subnet $\mathbf{a}^{(t)}$ is categorized into the w' -th supernet when

$$w' = \arg \min_{w \in [W]} \mathcal{L}(\boldsymbol{\theta}^{(t,w)}, \mathbf{z}, \mathbf{a}, \varepsilon_a).$$



Relationship with the adversarial bandit problems

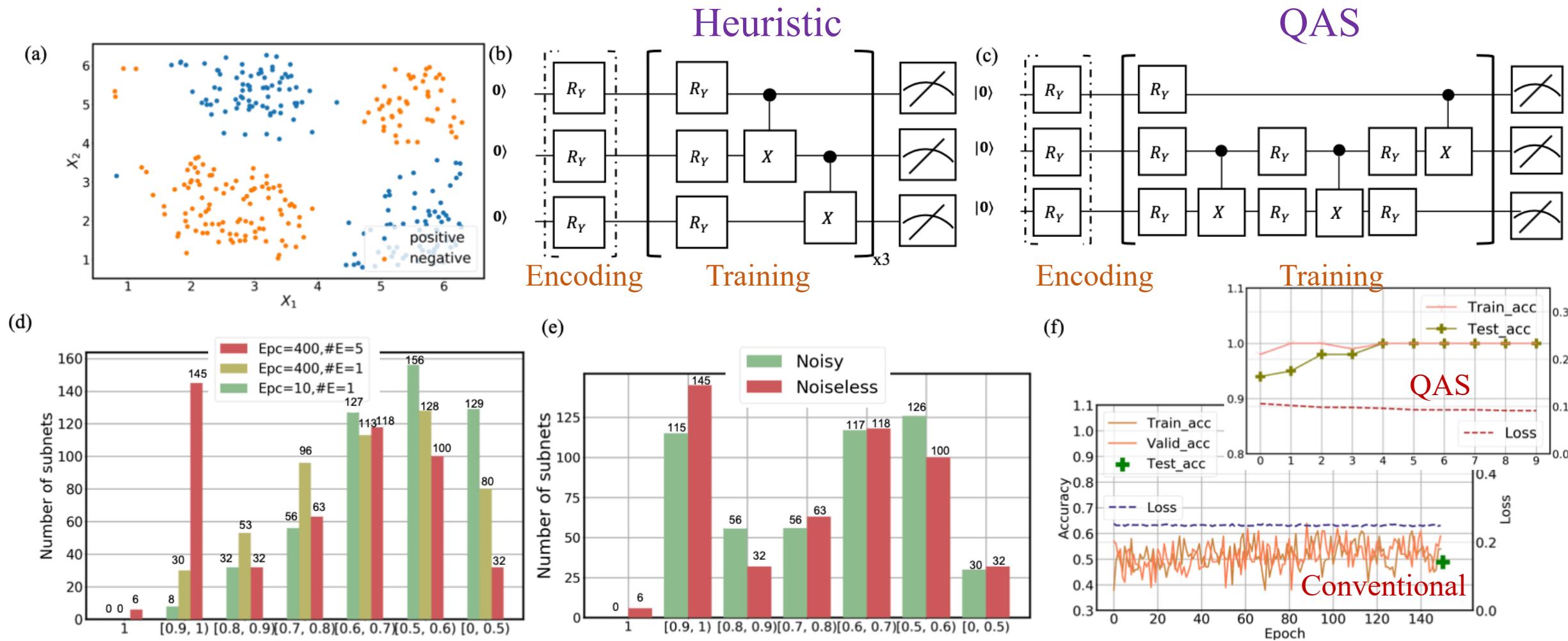
Theorem 1. Let W and T be the number of supernets and iterations, respectively. Suppose that the subnet $\mathbf{a}^{(t)}$ is assigned to the $I_w^{(t)}$ -th supernet $\mathcal{A}^{(I_w^{(t)})}$ with $I_w^{(t)} \in [W]$ at the t -th iteration, where the corresponding objective function is $\mathcal{L}(\boldsymbol{\theta}^{(t, I_w^{(t)})}, \mathbf{z}, \mathbf{a}, \varepsilon_a) \in [0, 1]$. Define the regret as

$$R_T = \sum_{t=1}^T \mathcal{L}(\boldsymbol{\theta}^{(t, I_w^{(t)})}, \mathbf{z}, \mathbf{a}, \varepsilon_a) - \min_{\{w_t\}_{t=1}^T} \sum_{t=1}^T \mathcal{L}(\boldsymbol{\theta}^{(t, I_w^{(t)})}, \mathbf{z}, \mathbf{a}, \varepsilon_a).$$

The method used in QAS to determine $\{I_w^{(t)}\}$ promises the regret $R_T \leq 0$, while the regret for the best bandit algorithms is lower bounded by $R_T = \Omega(T)$.

Summary: The proposed strategy outperforms all bandit learning algorithms in the measure of the regret.

Experiment results for classification tasks



Experiment results for the ground state energy estimation

JDT 京东科技

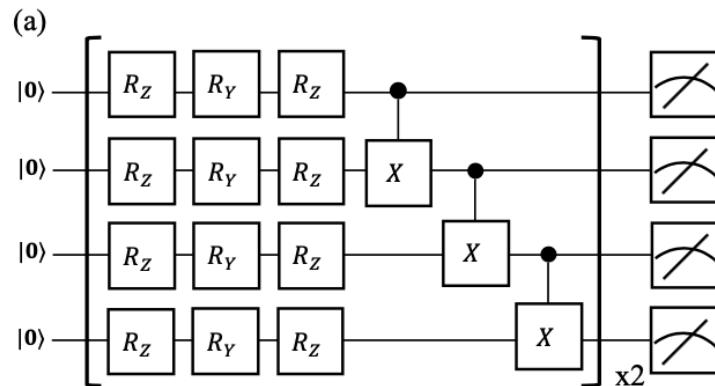
The molecular hydrogen Hamiltonian is formulated as

$$H_h = g + \sum_{i=1}^3 g_i Z_i + \sum_{i=1, k=1, i < k}^3 g_{i,k} Z_i Z_k + g_a Y_0 X_1 X_2 Y_3 + g_b Y_0 Y_1 X_2 X_3 + g_c X_0 X_1 Y_2 Y_3 + g_d X_0 Y_1 Y_2 X_3,$$

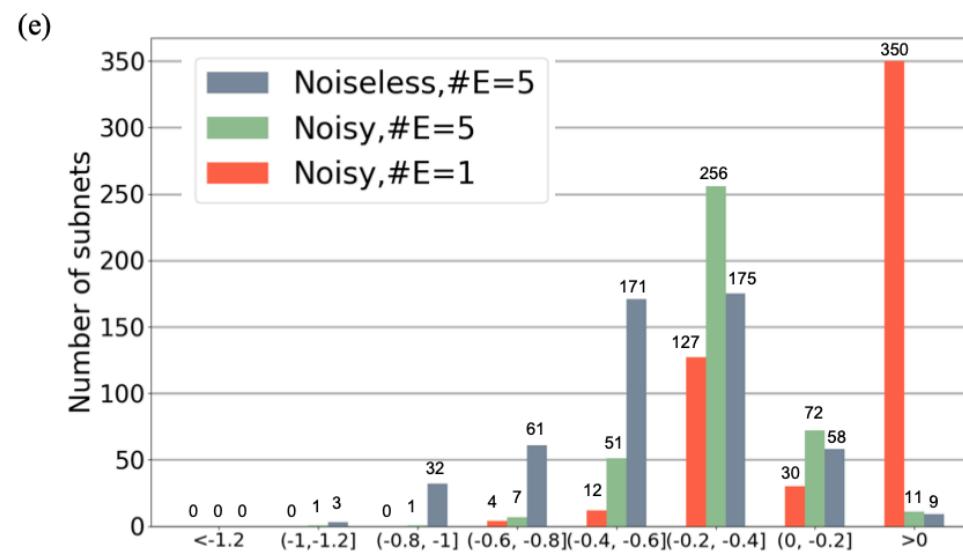
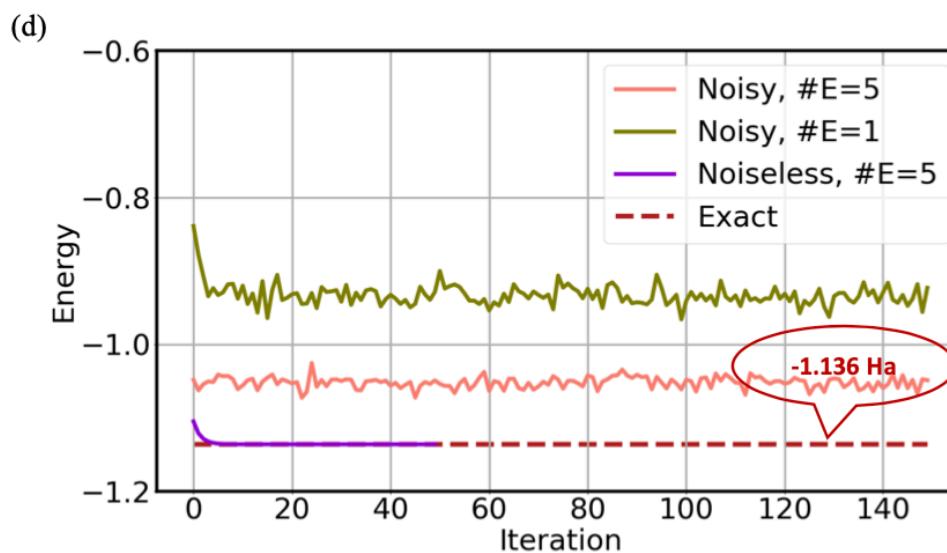
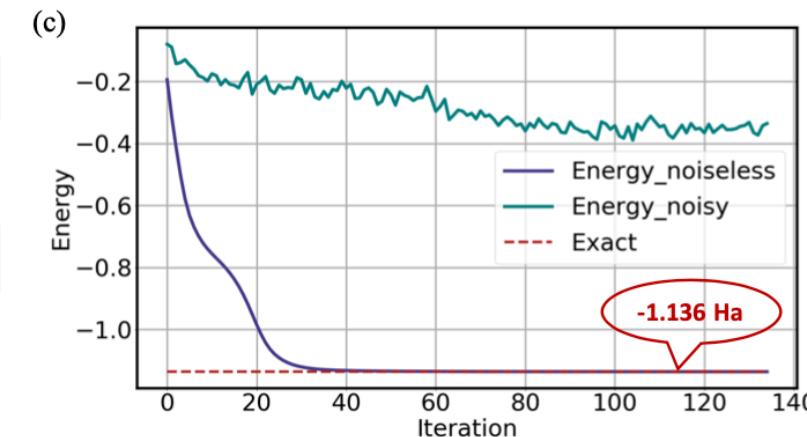
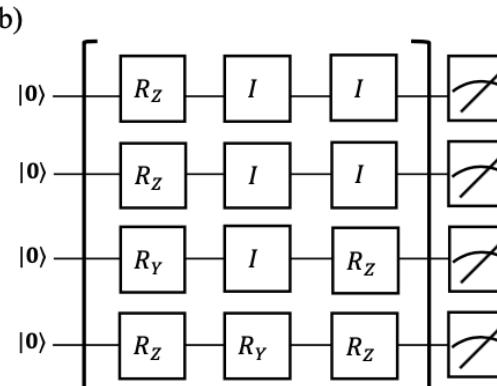
where $\{X_i, Y_i, Z_i\}$ denote the Pauli matrices acting on the i -th qubit and the real scalars g with or without subscripts are efficiently computable functions of the hydrogen-hydrogen bond length.

Experiment results for the ground state energy estimation

Heuristic



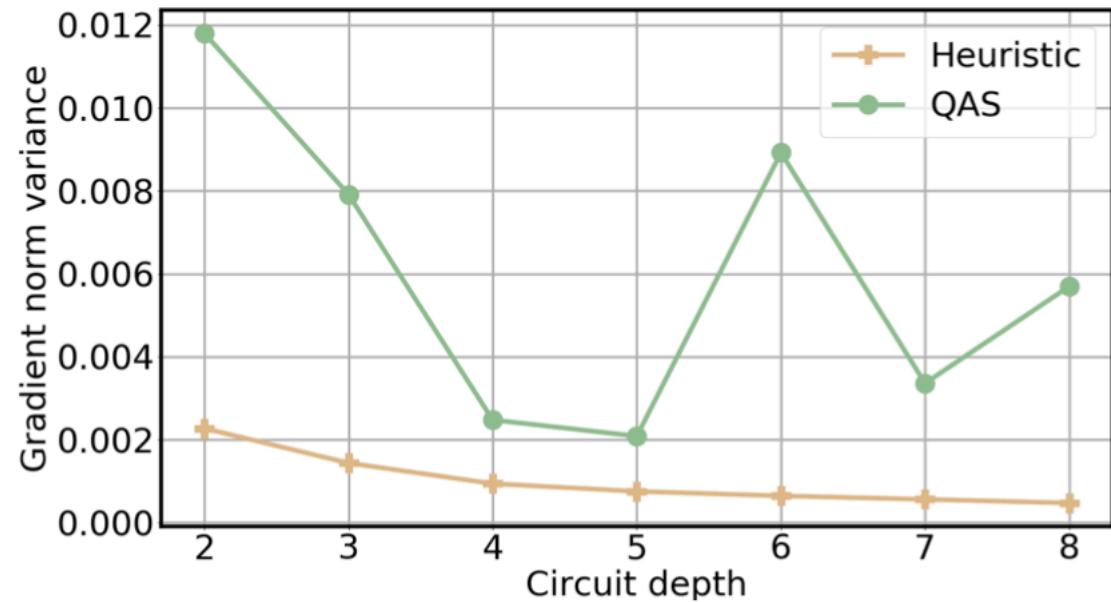
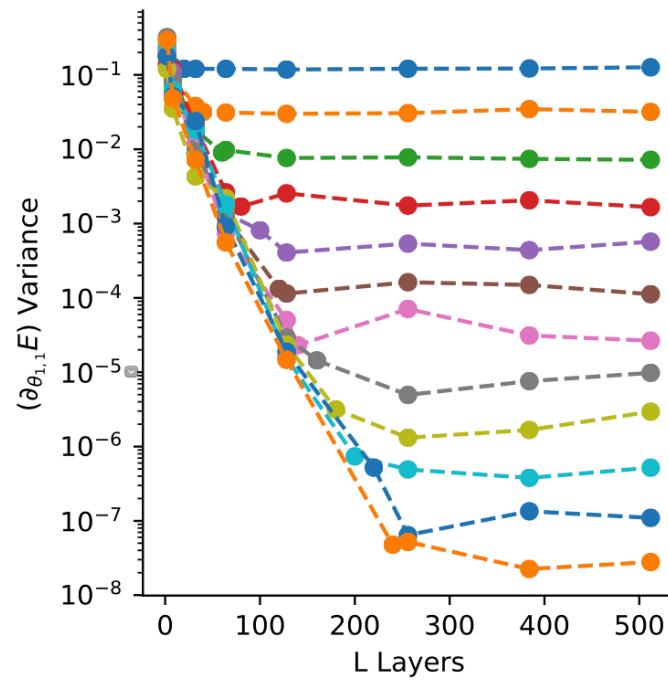
QAS



Experiment results for the ground state energy estimation

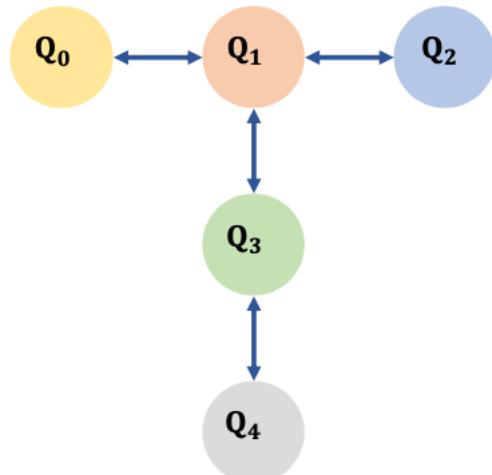
JDT 京东科技

QAS contributes to the alleviation of barren plateaus.



Experiment results for the ground state energy estimation

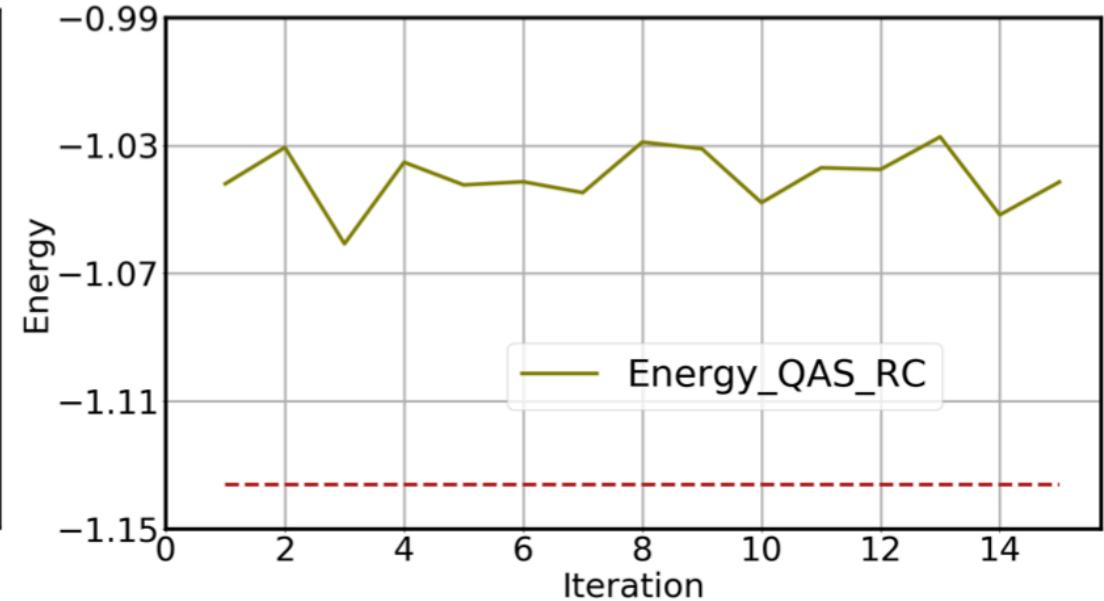
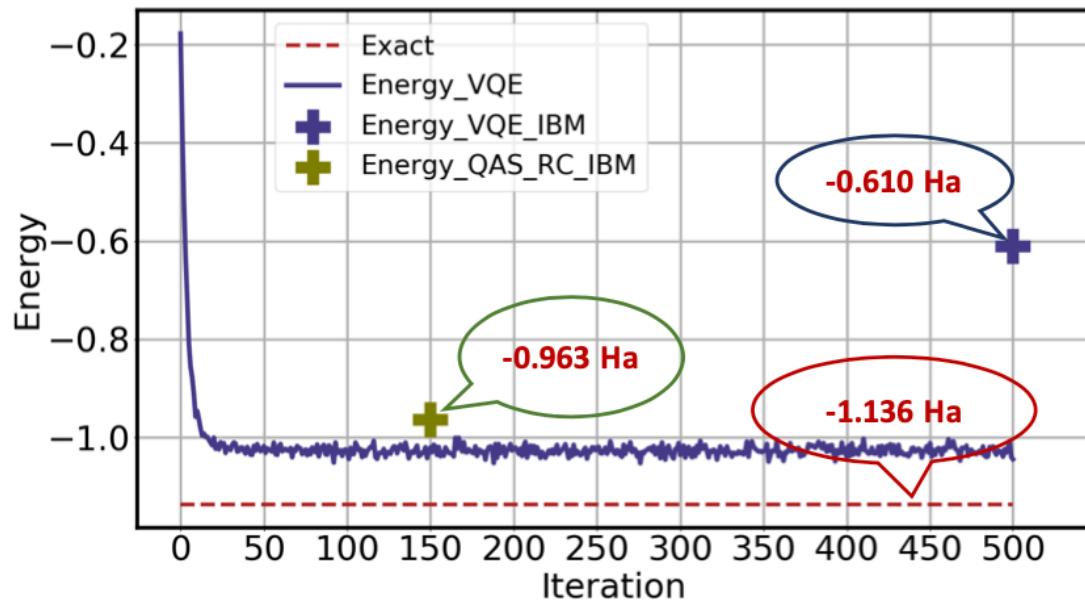
We last carry out QAS and the conventional VQE on IBM's 5-qubit quantum machine, i.e., 'Ibmq_ourense', to accomplish the ground state energy estimation of H_h .



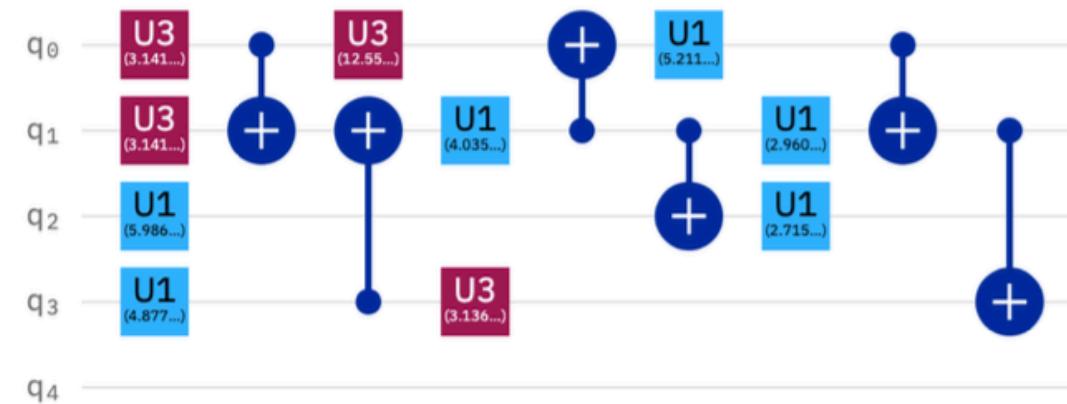
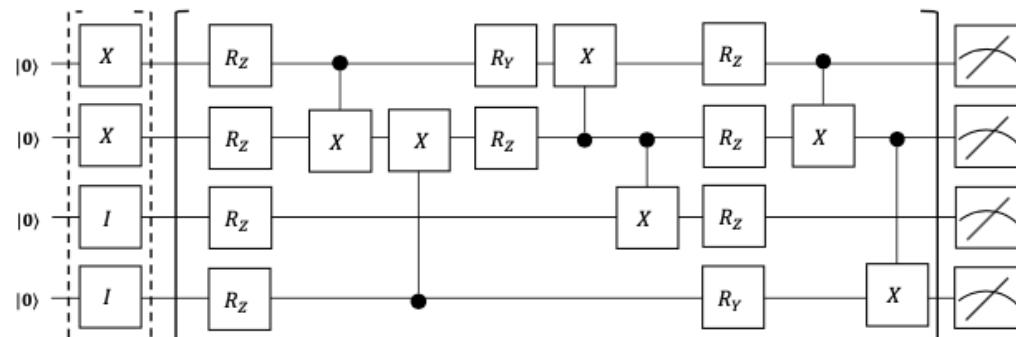
| Qubit | T1(μs) | T2(μs) | Readout error | Single-qubit U2 error gate | CNOT error rate |
|-------|---------------|---------------|---------------|----------------------------|--|
| Q0 | 75.75 | 50.81 | 1.65E-2 | 5.22E-4 | cx0-1: 9.55E-3 |
| Q1 | 78.47 | 27.56 | 2.38E-2 | 4.14E-4 | cx1-0: 9.55E-3 cx1-2: 9.44E-3 cx1-3: 1.25E-2 |
| Q2 | 101.51 | 107.00 | 1.57E-2 | 1.83E-4 | cx2-1: 9.44E-3 |
| Q3 | 79.54 | 78.38 | 3.95E-2 | 4.30E-4 | cx3-1: 1.25E-2 cx3-4: 8.34E-3 |
| Q4 | 74.27 | 30.00 | 4.74E-2 | 4.20E-4 | cx4-3: 8.34E-3 |

Experiment results for the ground state energy estimation

Experiment results of the ground state energy estimation of VQE and QAS.

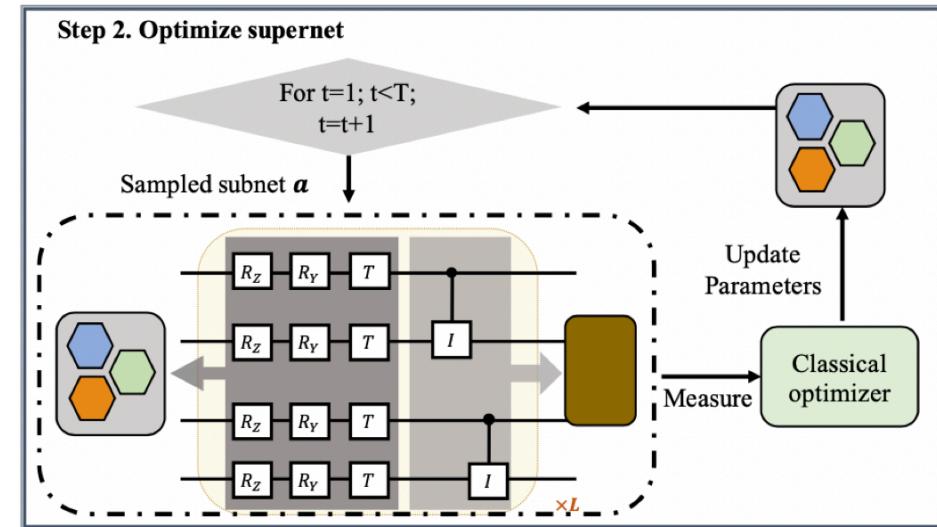


The searched circuit architecture by QAS



Improvements

- Improving the ranking stage of QAS, i.e., replacing the uniformly sampling by evolutionary algorithms or reinforcement learning algorithms.
- Employing the adversarial bandit learning algorithms to train multiple supernets. Benefits: reduce the runtime cost. Disadvantages: inducing a relatively large regret bound.

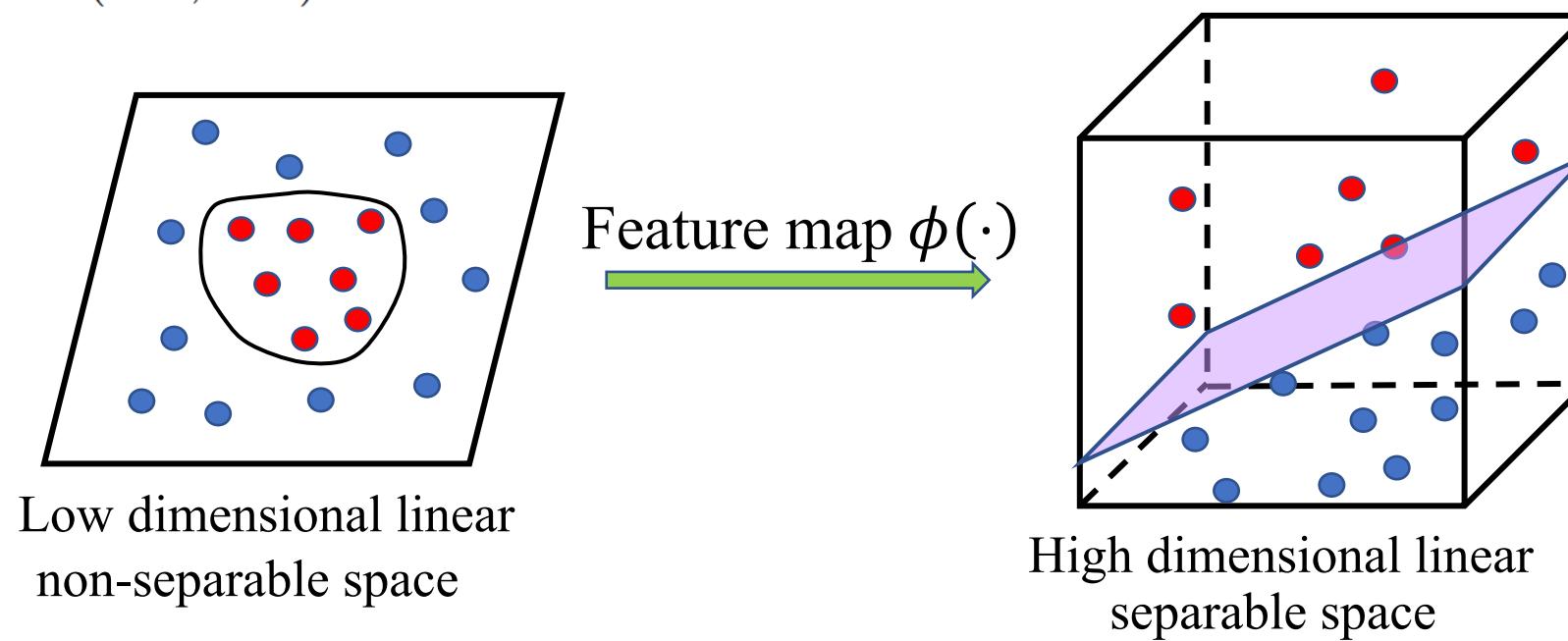


III. The power of quantum kernels in the NISQ era

Kernel methods

Core idea:

- Mapping the given input $x^{(i)} \in \mathbb{R}^d$ into a high-dimensional feature space, i.e., $\phi(\cdot): \mathbb{R}^d \rightarrow \mathbb{R}^q$ ($q \gg d$).
- Constructing a kernel function $\kappa(x^{(i)}, x^{(j)}) = \langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle$ and the corresponding kernel matrix $K_{ij} = \kappa(x^{(i)}, x^{(j)})$



Typical kernels: the radial basis function kernel, Gaussian kernel, and polynomial kernel.

Kernel methods

Core idea:

For efficient optimization of w , we consider minimization of the following loss function

$$\min_w \lambda \langle w, w \rangle + \sum_{i=1}^N \left(\langle w, \phi(x_i) \rangle - \text{Tr}(U^\dagger O U \rho(x_i)) \right)^2,$$

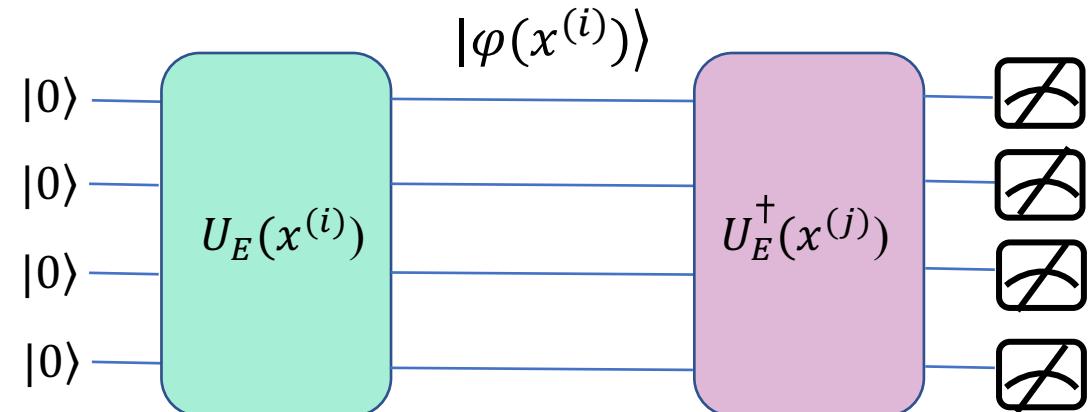
$\text{h}(x_i)$ y_i

The optimal w can be written down explicitly as

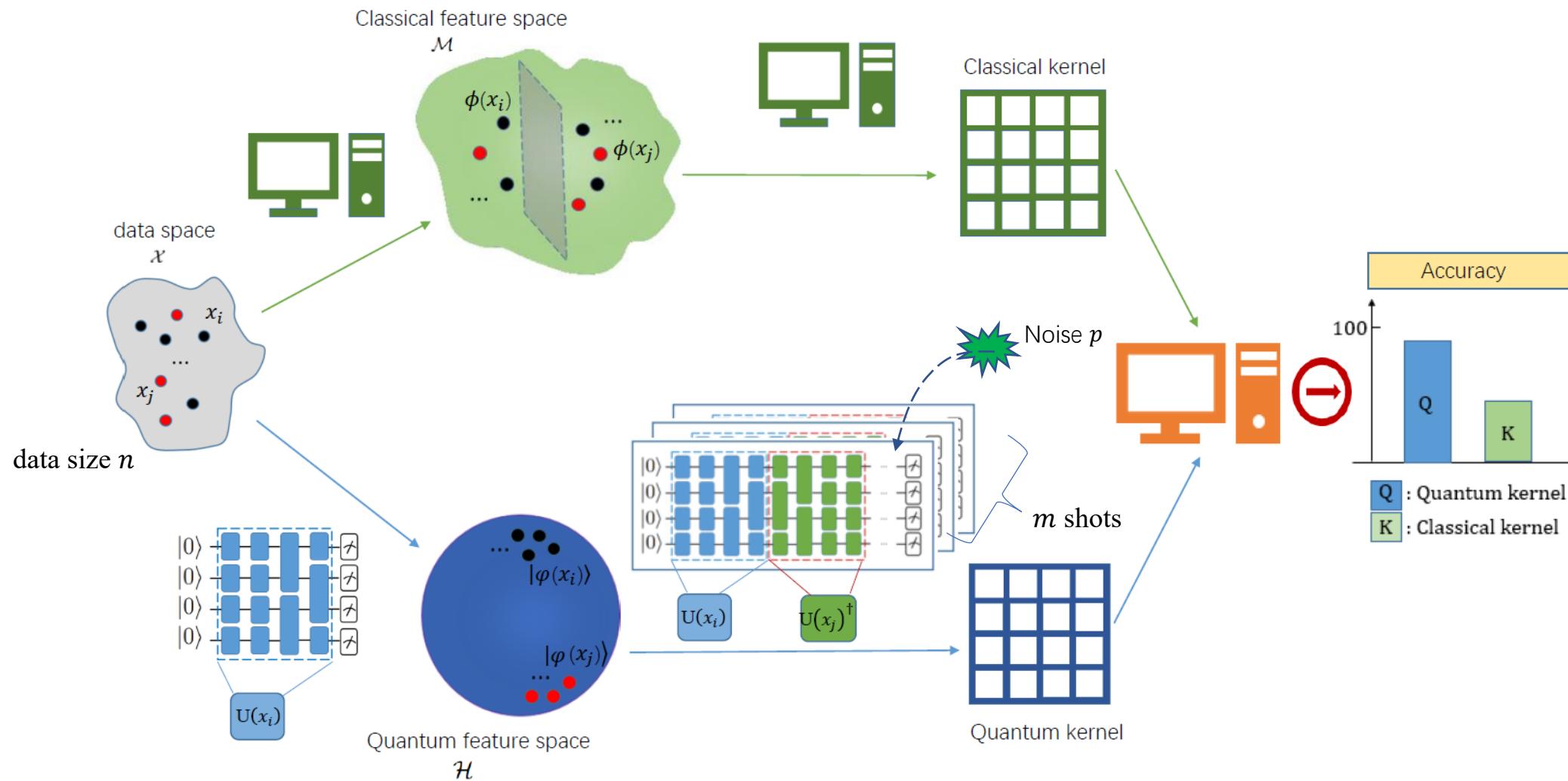
$$w = \sum_{i=1}^N \sum_{j=1}^N \phi(x_i) ((K + \lambda I)^{-1})_{ij} \text{Tr}(U^\dagger O U \rho(x_j)).$$

Quantum kernel methods

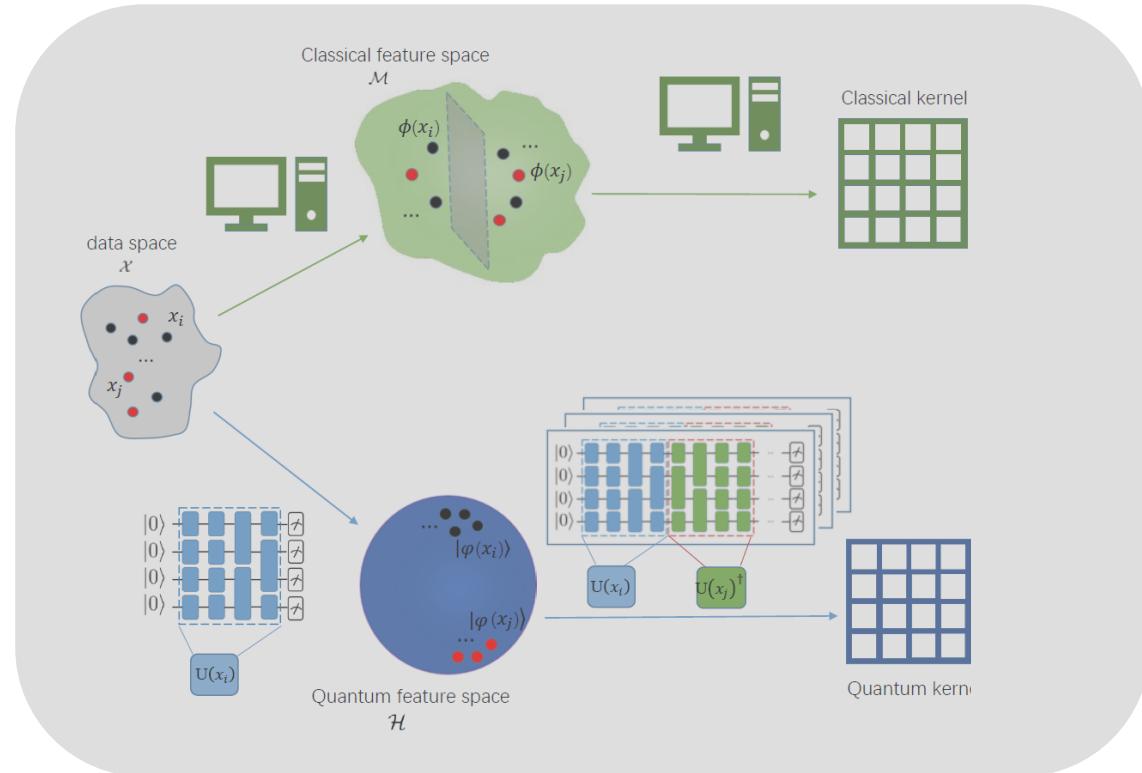
- Replacing classical feature map $\phi(\cdot)$ with a variational quantum circuit $U_E(\cdot)$.
- Mapping the give classical data $x^{(i)} \in \mathbb{R}^d$ to exponentially large feature space, i.e., $|\phi(x^{(i)})\rangle = U_E(x^{(i)})|0^d\rangle \in \mathbb{R}^{2^d}$. Denote N as the number of qubits, then $N = d$.
- The result of kernel function $\kappa(\cdot, \cdot)$ coincides with applying measurements on the prepared. quantum states, then the (i, j) -th element in kernel matrix is $\mathbf{W}_{ij} = |\langle \varphi(\mathbf{x}^{(j)}) | \varphi(\mathbf{x}^{(i)}) \rangle|^2$



Overview



The power of quantum kernels



How quantum kernel methods outperform classical kernel methods in predicting the unseen data ?

Problem setup.

Each data pair $(x^{(i)}, y^{(i)})$ is constructed by $y^{(i)} = f(x^{(i)}) = \text{Tr}(O U(\boldsymbol{\theta}^*) \rho(x^{(i)}) U(\boldsymbol{\theta}^*)^\dagger)$, where $U(\boldsymbol{\theta}^*)$ a specified unitary evolution applied on the encoded quantum example $\rho(x^{(i)})$, and O is a certain measurement operator.

Given N examples, the goal of a learning model is to infer a map $h(\cdot)$ such that $y^{(i)} = h(x^{(i)})$.

Results.

Huang et al. [arXiv:2011.01938] present a generalization error bound under **the ideal scenario**, i.e., given n examples, the generalization error of classical kernels follows

$$\mathbb{E}_{x \sim \mathcal{D}} |h(x) - f(x)| \leq c \sqrt{\frac{s_K(n)}{n}}$$

where $s_K(n)$ refers to the model complexity with

$$s_K(n) = \sum_{i=1}^n \sum_{j=1}^n (K^{-1})_{ij} \text{Tr}(O^U \rho(x_i)) \text{Tr}(O^U \rho(x_j)).$$

The power of data in quantum Kernels

JDT 京东科技

$$\mathbb{E}_{x \sim \mathcal{D}} |h(x) - f(x)| \leq c \sqrt{\frac{s_K(n)}{n}}$$
$$s_K(n) = \sum_{i=1}^n \sum_{j=1}^n (K^{-1})_{ij} \text{Tr}(O^U \rho(x_i)) \text{Tr}(O^U \rho(x_j)).$$

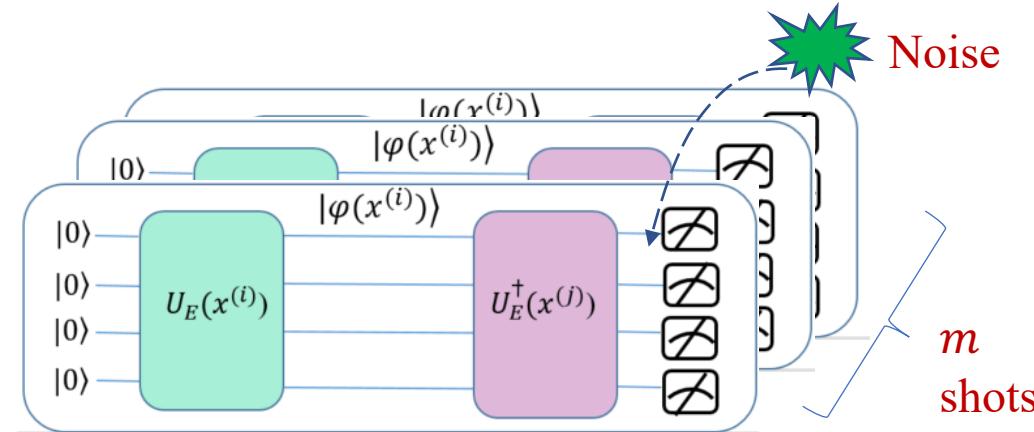
Two insights

- Given the same sample size, a small $\frac{Y^\top W^{-1} Y}{Y^\top K^{-1} Y}$ leads to a large quantum advantage over the classical kernel methods.
- Under the same kernel method, a large n leads to a better prediction performance.

How about power of quantum kernels in the NISQ era?

JDT 京东科技

In the NISQ era, quantum kernels are affected by **the system noise** and **finite shots**.



- Consider the depolarization channel $\mathcal{N}_p(\rho(x^{(i)})) = (1 - p)\rho(x^{(i)}) + \frac{p\mathbb{I}_{2^N}}{2^N}$, where p is the depolarization rate. Then the quantum kernel can be written as

$$\widetilde{\mathbf{W}} = (1 - p) \mathbf{W} + \frac{p\mathbb{I}_{2^N}}{2^N}$$

- Consider the finite measurements, the estimated element of quantum kernels yields

$$\widehat{\mathbf{W}}_{ij} = \frac{1}{m} \sum_{k=1}^m V_k, \quad V_k \sim \text{Ber}(\widetilde{\mathbf{W}}_{ij})$$

The power of quantum kernels under the NISQ setting

JDT 京东科技

Under the NISQ setting, we give a generalization error bound for the estimate kernel \widehat{W} .

Theorem 3. Let the size of training dataset be n and the number of measurements is m . Define $Y = [y_1, \dots, y_n]^\top$ as the label vector and $c_W = \|W^{-1}\|_2$. Suppose the system noise is modeled by \mathcal{N}_p . With probability at least $1 - \delta$, the noisy quantum kernel \widehat{W} can be used to infer a hypothesis $h(\mathbf{x})$ with generalization error

$$\mathbb{E}_{\mathbf{x}, \widehat{W}} |h(\mathbf{x}) - f(\mathbf{x})| \leq \tilde{O} \left(\sqrt{\frac{c_1}{n}} + \sqrt{\frac{1}{c_2} \frac{n}{\sqrt{m}}} \right) \quad (\text{D10})$$

where $c_1 = Y^\top W^{-1} Y$ and $c_2 = \max(c_W^{-2} ((\frac{1}{2} \log(\frac{4n^2}{\delta}))^{\frac{1}{2}} + m^{\frac{1}{2}} p (1 + \frac{1}{2^{N+1}}))^{-1} - \frac{n}{\sqrt{m}} c_W^{-1}, 0)$.

An additional term $\sqrt{n/(c_2 \sqrt{m})}$ compared to the bound of Huang et al. has following implications:

- n/\sqrt{m} suggests that an increased number of data n will result in a higher generalization error.
- The performance of quantum kernels in the NISQ era are heavily depends on the number of measurements.
- The term c_2 indicates the negative role of the system noise. Moreover, the generalization error bound would be infinite once $p > 1/(nc_W(1 + 1/2^{N+1}))$.

The power of quantum kernels under the NISQ setting

JDT 京东科技

Under the NISQ setting, we give a generalization error bound for the estimate kernel \widehat{W} .

$$\mathbb{E}_{\mathbf{x}, \widehat{W}} |h(\mathbf{x}) - f(\mathbf{x})| \leq \tilde{O} \left(\sqrt{\frac{c_1}{n}} + \sqrt{\frac{1}{c_2} \frac{n}{\sqrt{m}}} \right) \quad (\text{D10})$$

where $c_1 = Y^\top W^{-1} Y$ and $c_2 = \max(c_W^{-2}((\frac{1}{2} \log(\frac{4n^2}{\delta}))^{\frac{1}{2}} + m^{\frac{1}{2}} p (1 + \frac{1}{2^{N+1}}))^{-1} - \frac{n}{\sqrt{m}} c_W^{-1}, 0)$.

An additional term $\sqrt{n/(c_2\sqrt{m})}$ compared to the bound of Huang et al. has following implications:

- n/\sqrt{m} suggests that an increased number of data n will result in a higher generalization error;
- The performance of noisy quantum kernels heavily depends on shots m ;
- The term c_2 indicates the negative role of the system noise. Moreover, the generalization error bound would be infinite once $p > 1/(nc_W(1 + 1/2^{N+1}))$.

Proof Sketches of Theorem 1

Following the result of Huang et al., with probability $1 - \delta/2$, the generalization error of \widehat{W} yields

$$\mathbb{E}_{\mathbf{x}, \widehat{W}} |f(\mathbf{x}) - h(\mathbf{x})| \leq \tilde{O} \left(\sqrt{\frac{Y^\top \widehat{W}^{-1} Y}{n}} \right)$$

$$\text{Note that } \sqrt{Y^\top \widehat{W}^{-1} Y} = \sqrt{Y^\top (\widehat{W}^{-1} - W^{-1}) Y + Y^\top W^{-1} Y} \leq \sqrt{\|\widehat{W}^{-1} - W^{-1}\|_2 Y^\top Y} + \sqrt{Y^\top W^{-1} Y} \quad (3)$$

and the connection of $\|\widehat{W}^{-1} - W^{-1}\|_2$ and $\|W - \widehat{W}\|_2$, that is

$$\|W^{-1} - \widehat{W}^{-1}\|_2 \leq \frac{c_W^2 \|W - \widehat{W}\|_2}{1 - c_W \|W - \widehat{W}\|_2}. \quad (c_W = \|W^{-1}\|_2) \quad (4)$$

To achieve $\|W^{-1} - \widehat{W}^{-1}\|_2 \leq \delta$, it is sufficient to show $\|W - \widehat{W}\|_2 \leq \frac{\delta}{c_W(\delta + c_W)}$.

Meanwhile, utilizing the Chernoff-Hoeffding bound and the equation $\widetilde{W} = (1 - p)W + \frac{p\mathbb{I}_{2^N}}{2^N}$ yields

$$\Pr \left(|W_{ij} - \widehat{W}_{ij}| \geq p \left(1 + \frac{1}{2^{N+1}} \right) + \frac{\delta}{2} \right) \leq 2 \exp \left(\frac{-\delta^2 m}{2} \right). \quad (5)$$

Proof Sketches of Theorem 1

Then a similar concentration inequality in term of the norm $\|W - \widehat{W}\|_2$ can be written as

$$\begin{aligned}
 \Pr(\|W - \widehat{W}\|_2 \geq \delta') &\leq \Pr(\|W - \widehat{W}\|_F \geq \delta') && (\|W - \widehat{W}\|_2 \leq \|W - \widehat{W}\|_F) \\
 &= \Pr\left(\sum_{i=1}^n \sum_{j=1}^n |W_{ij} - \widehat{W}_{ij}|^2 \geq \delta'^2\right) \\
 &\leq \Pr\left(\bigcup_{i=1}^n \bigcup_{j=1}^n |W_{ij} - \widehat{W}_{ij}|^2 \geq \frac{\delta'^2}{n^2}\right) && \text{(union bound)} \\
 &\leq 2n^2 \exp\left(-2\left(\frac{\delta'}{n} - p\left(1 + \frac{1}{2^{N+1}}\right)\right)^2 m\right) && \text{(using Eqn. (4) and requiring } \delta' > np\left(1 + \frac{1}{2^{N+1}}\right)\text{)}
 \end{aligned}$$

Combining this inequality with Eqn. (4), we have

$$\Pr(\|W^{-1} - \widehat{W}^{-1}\|_2 \geq \delta) \leq \Pr\left(\|W - \widehat{W}\|_2 \geq \frac{\delta}{c_W(\delta + c_W)}\right) \leq 2n^2 \exp\left(-2\left(\frac{\delta}{nc_W(\delta + c_W)} - p\left(1 + \frac{1}{2^{N+1}}\right)\right)^2 m\right). \quad (6)$$

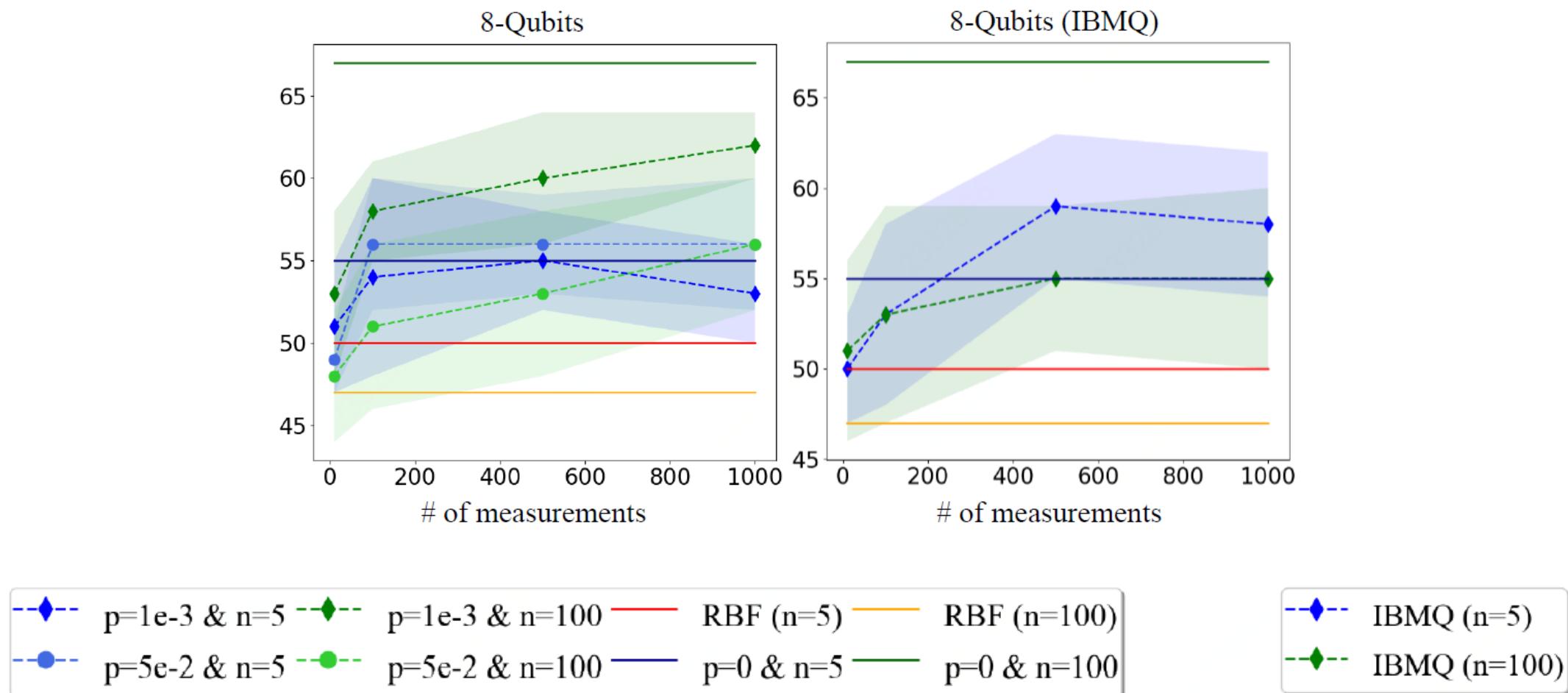
Setting the right term in Eqn. (6) as $\delta'/2$, we can get that with probability at least $1 - \delta'/2$, (needs to meet non-negative)

$$\|W^{-1} - \widehat{W}^{-1}\|_2 \leq \max\left(c_W^{-2} \left(\left(\frac{1}{2} \log\left(\frac{4n^2}{\delta}\right)\right)^{\frac{1}{2}} + m^{\frac{1}{2}} p\left(1 + \frac{1}{2^{N+1}}\right)\right)^{-1} - \frac{n}{\sqrt{m}} c_W^{-1}, 0\right).$$

Note that this bound can meet all needed conditions. And this completes the proof.

Numerical evidence

The numerical simulations conducted on the Fashion-Mnist dataset accord with our theoretical results.



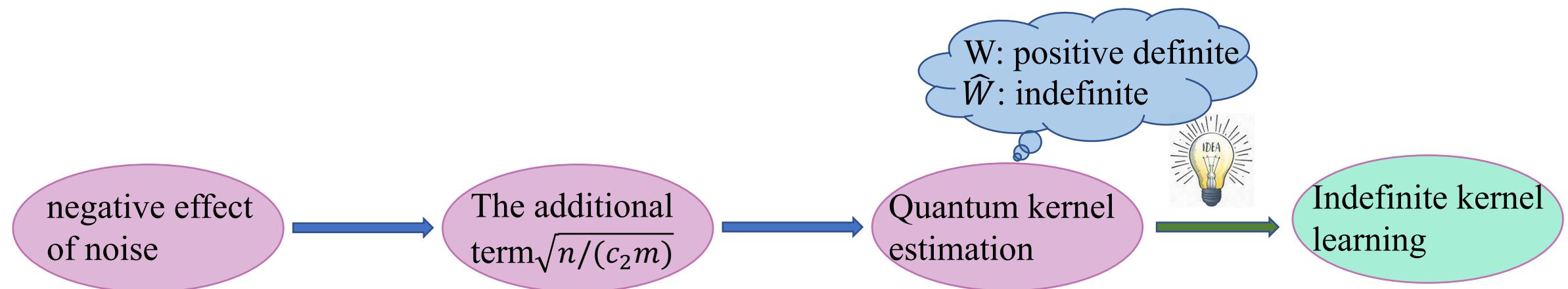
Enhance Performance of Noise Quantum Kernels

JDT 京东科技

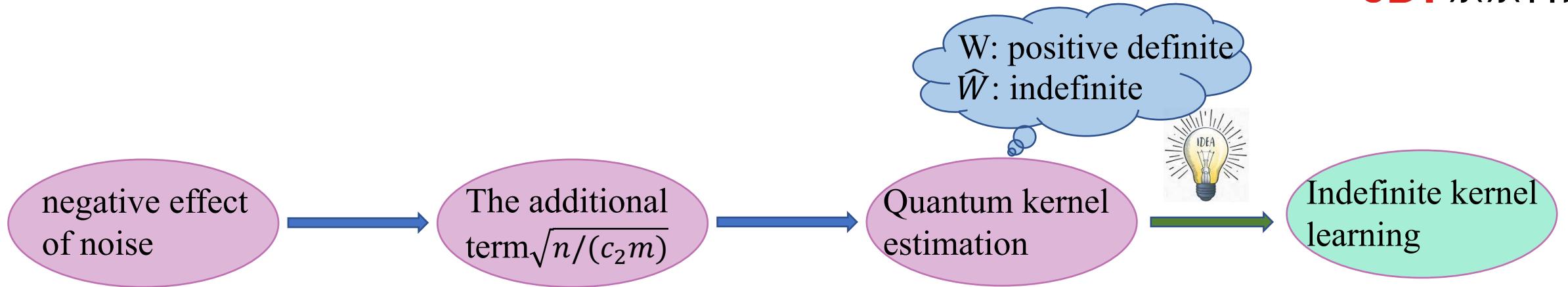
how to preserve the superiority of quantum kernels carried on NISQ machines?

- Slimming the size of training data;
- Suppressing the effects of noise (**HOW?**)

Tracing back the source of such negative effect can get that



Enhance Performance of Noise Quantum Kernels



Denote the eigenvalues of \widehat{W} as $\widehat{\lambda}_i$ ($\widehat{\lambda}_1 \geq \dots \geq \widehat{\lambda}_r \geq 0 \geq \dots \geq \widehat{\lambda}_n$) and the eigenvectors as \mathbf{u}_i .

Three advanced spectral transformation techniques:

Clip: clipping all negative eigenvalues to zero, i.e., $\widehat{W}_c = \sum_{i=1}^r \widehat{\lambda}_i \mathbf{u}_i \mathbf{u}_i^\top$.

Flip: flipping the sign of negative eigenvalues, i.e., $\widehat{W}_f = \sum_{i=1}^r \widehat{\lambda}_i \mathbf{u}_i \mathbf{u}_i^\top - \sum_{i=1}^r \widehat{\lambda}_i \mathbf{u}_i \mathbf{u}_i^\top$.

Shift: shifting all eigenvalues by a positive constant, i.e., $\widehat{W}_s = \sum_{i=1}^n (\widehat{\lambda}_i - \widehat{\lambda}_n) \mathbf{u}_i \mathbf{u}_i^\top$.

Note that our goal is to shrink the distance between noiseless and noisy quantum kernels.

Enhance Performance of Noise Quantum Kernels

JDT 京东科技

Lemma 1. Let W and \widehat{W} be the ideal and noisy quantum kernel, respectively. Applying the spectral transformation techniques to \widehat{W} , the obtained kernel $\widehat{W}_\diamond \in \{\widehat{W}_c, \widehat{W}_f, \widehat{W}_s\}$ yields

$$\|W - \widehat{W}_\diamond\|_F \leq \|W - \widehat{W}\|_F, \quad (7)$$

where $\|\cdot\|_F$ refers to the Frobenius norm.

Proof sketches: We give a brief proof for the clipping case and the other cases are similar.

Denote $W = \sum_{i=1}^n \lambda_i v_i v_i^\top$, where λ_i refer to the eigenvalues and v_i are the corresponding eigenvectors.

Supported by the definition of the Frobenius norm, an equivalent of achieving Eqn. (6) is

$$\text{Tr} \left((W - \widehat{W}_c)^2 \right) \leq \text{Tr} \left((W - \widehat{W})^2 \right) \longleftrightarrow \text{Tr} \left(\widehat{W}_c^2 \right) - 2 \text{Tr} \left(W \widehat{W}_c \right) \leq \text{Tr} \left(\widehat{W}^2 \right) - 2 \text{Tr} \left(W \widehat{W} \right)$$

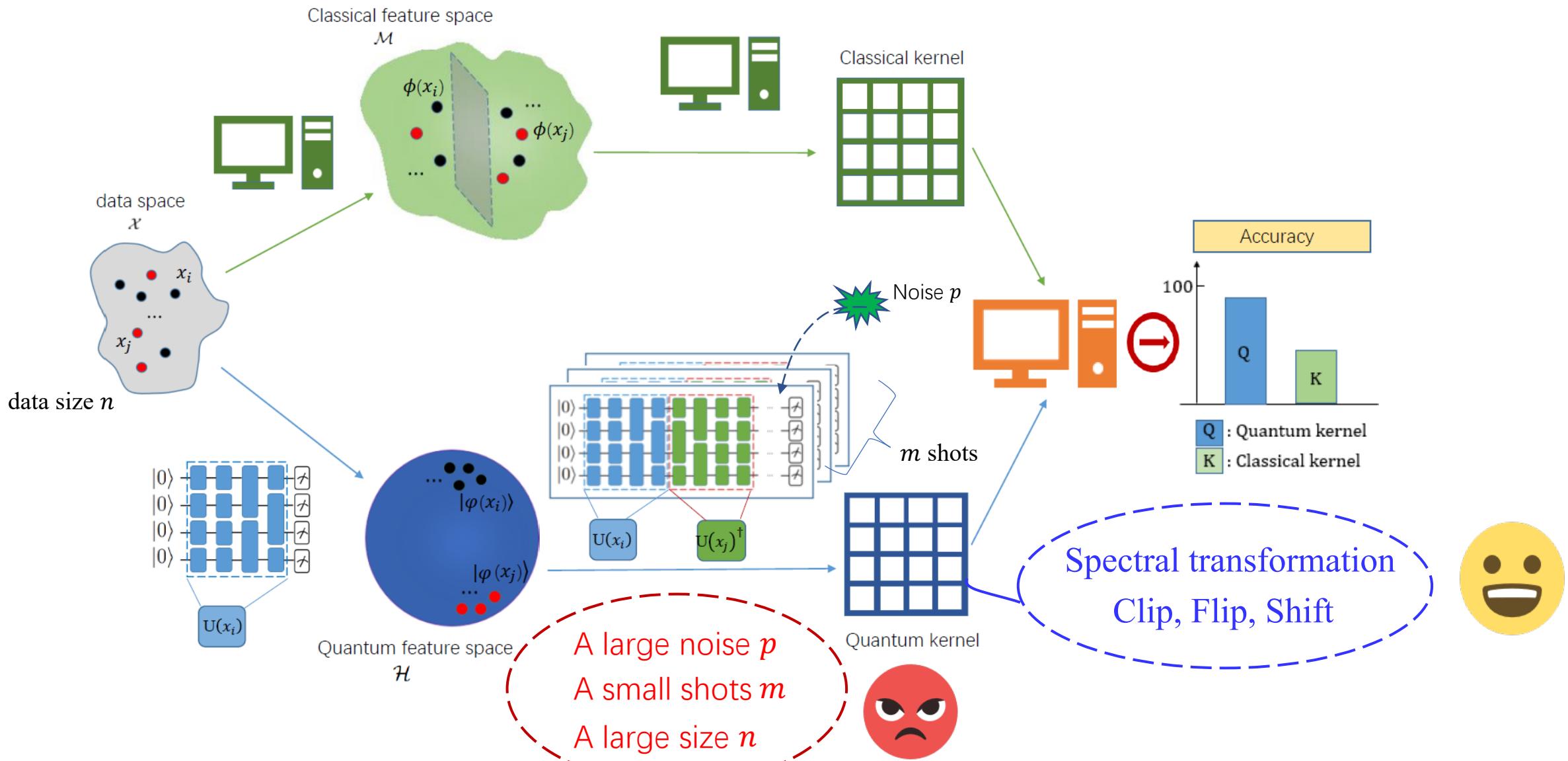
Considering $\text{Tr}(\widehat{W}^2)$, $\text{Tr}(\widehat{W}^2) = \sum_{i=1}^n \widehat{\lambda}_i^2 \geq \sum_{i=1}^r \widehat{\lambda}_i^2 = \text{Tr}(\widehat{W}_c^2)$. (8)

Considering $\text{Tr}(W \widehat{W})$,

$$\text{Tr}(W \widehat{W}) = \text{Tr} \left(\left(\sum_{i=1}^n \lambda_i v_i v_i^\top \right) \left(\sum_{j=1}^r \widehat{\lambda}_j u_j u_j^\top + \sum_{j=r+1}^n \widehat{\lambda}_j u_j u_j^\top \right) \right) = \text{Tr} \left(W \widehat{W}_c \right) + \sum_{i=1}^n \sum_{j=r+1}^n \lambda_i \widehat{\lambda}_j (u_j^\top v_i)^2 \leq \text{Tr} \left(W \widehat{W}_c \right). \quad (9)$$

In conjunction with Eqn. (8) and (9), we can achieve Eqn (7).

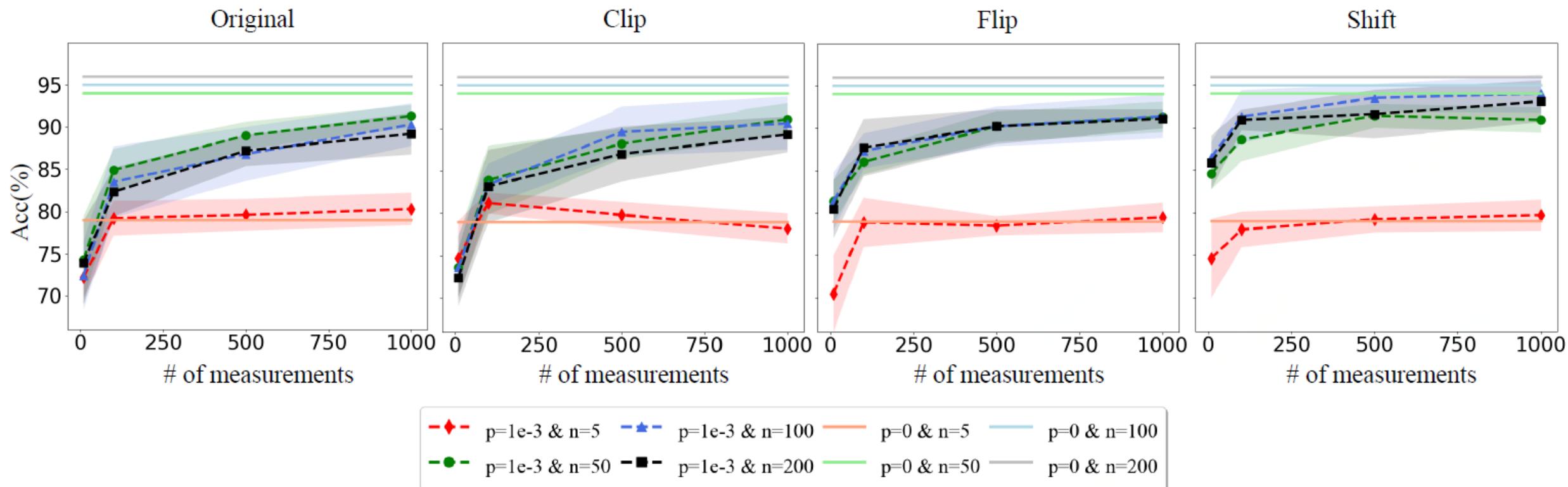
The power of noisy quantum kernels



Enhance Performance of Noise Quantum Kernels

JDT 京东科技

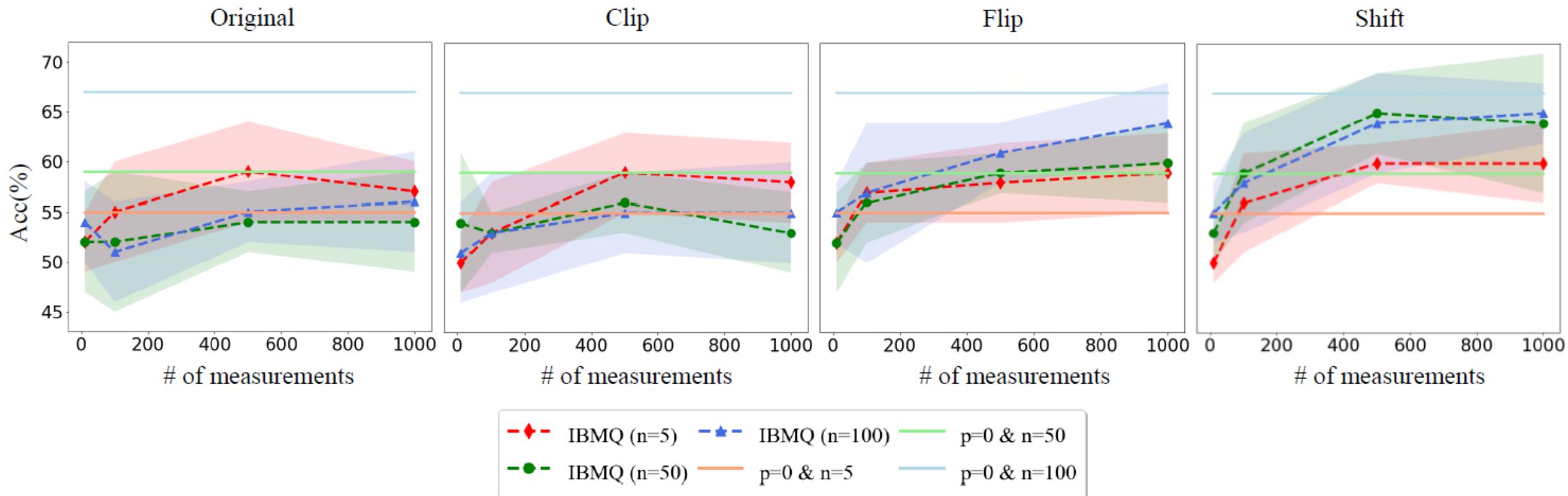
Numerical simulations achieved with the depolarization noise model



Enhance Performance of Noise Quantum Kernels

JDT 京东科技

Numerical simulations achieved on IBMQ-Melbourne's noisy settings.



Thank You!