

# Supplemental Material: Representation Learning via Quantum Neural Tangent Kernels

Junyu Liu,<sup>1,2,3,\*</sup> Francesco Tacchino,<sup>4,†</sup> Jennifer R. Glick,<sup>5,‡</sup> Liang Jiang,<sup>1,2,§</sup> and Antonio Mezzacapo<sup>5,¶</sup>

<sup>1</sup>*Pritzker School of Molecular Engineering, The University of Chicago, Chicago, IL 60637, USA*

<sup>2</sup>*Chicago Quantum Exchange, Chicago, IL 60637, USA*

<sup>3</sup>*Kadanoff Center for Theoretical Physics, The University of Chicago, Chicago, IL 60637, USA*

<sup>4</sup>*IBM Quantum, IBM Research, Zurich, 8803 Rüschlikon, Switzerland*

<sup>5</sup>*IBM Quantum, IBM T. J. Watson Research Center, Yorktown Heights, New York 10598, USA*

(Dated: July 13, 2022)

∗: Corresponding author.

## CONTENTS

I. Details on the linear and the quadratic models	1
A. Linear model	1
B. Quadratic model	3
1. The residual training error	5
2. The asymptotic regime	7
II. Details on the quantum optimization	8
A. General setup	8
B. Frozen QNTK	9
C. dQNTK	10
III. Details on the quantum machine learning: Hermitian operator expectation value evaluation	11
A. General setup	11
B. No representation learning	13
C. Representation learning	14
IV. Details on the quantum machine learning: amplitude encoding	16
A. General setup	16
B. No representation learning	18
C. Representation learning	19
D. Reading the amplitude	23
V. Suppression of non-Gaussianity in the large-width limit	24
VI. On the LeCun parametrization and the NTK parametrization	26
VII. Simulation details and additional data	27
References	28

## I. DETAILS ON THE LINEAR AND THE QUADRATIC MODELS

### A. Linear model

Before we discuss the quantum model setup, we start by reviewing the linear and quadratic models in the classical setup. The linear and quadratic models are originally solved in the framework presented by [1, 2], and in this section, we will give a brief review.

---

∗ junyuliu@uchicago.edu

† fta@zurich.ibm.com

‡ jennifer.r.glick@ibm.com

§ liang.jiang@uchicago.edu

¶ mezzacapo@ibm.com

We define the model as

$$z_{i;\delta}(\theta) = \sum_j W_{ij} \phi_j(\mathbf{x}_\delta) . \quad (1)$$

Here  $\mathbf{x}_\delta$  is the data point  $\delta$  in the space  $\mathcal{D}$ , and  $W_{ij}$ s are weights and biases. We use the slack notation such that  $W_{0j} = b_j$  includes the bias. The model is called the linear model, which is linear in weights, but we wish to keep the feature map  $\phi$  in general. Moreover, we write  $\theta$  as a compact notation of the vectorized  $W$ .

We wish to optimize the following loss function

$$\mathcal{L}_{\mathcal{A}}(\theta) = \frac{1}{2} \sum_{i,\tilde{\alpha}} \left[ y_{i;\tilde{\alpha}} - \sum_j W_{ij} \phi_j(\mathbf{x}_{\tilde{\alpha}}) \right]^2 . \quad (2)$$

For the sample set  $\mathcal{A}$ . Again, we will use the notation  $\tilde{\alpha}$  such that  $\tilde{\alpha}$  is in the sample set  $\mathcal{A}$ . We define the kernel

$$k_{\delta_1 \delta_2} \equiv k(\mathbf{x}_{\delta_1}, \mathbf{x}_{\delta_2}) \equiv \sum_i \phi_i(\mathbf{x}_{\delta_1}) \phi_i(\mathbf{x}_{\delta_2}) = \sum_{a,b} \frac{dz_{i;\delta_1}}{dW_{ab}} \frac{dz_{i;\delta_2}}{dW_{ab}} . \quad (3)$$

Note that the right hand side does not contain a sum over  $i$ . It is equal for all  $i$ s. Moreover, We call  $\delta \in \mathcal{D}$  as the whole data set, while  $\mathcal{A} \subset \mathcal{D}$  is the sample set. We define  $\tilde{k}_{\tilde{\alpha}_1 \tilde{\alpha}_2}$  with tilde to indicate the submatrix. We also define [3]

$$\sum_{\tilde{\alpha}_2 \in \mathcal{A}} \tilde{k}^{\tilde{\alpha}_1 \tilde{\alpha}_2} \tilde{k}_{\tilde{\alpha}_2 \tilde{\alpha}_3} = \delta_{\tilde{\alpha}_1 \tilde{\alpha}_3}^{\tilde{\alpha}_1} . \quad (4)$$

namely, the upper index means the inverse. Similar to the main text, we consider the gradient descent algorithm,

$$d\theta_\mu = -\eta \frac{\partial \mathcal{L}_{\mathcal{A}}}{\partial \theta_\mu} . \quad (5)$$

The partial derivatives are computed as

$$\begin{aligned} \frac{\partial L_{\mathcal{A}}(W)}{\partial W_{ab}} &= - \sum_{\tilde{\alpha}, i, j} \delta_{ia} \delta_{jb} \phi_j(\mathbf{x}_{\tilde{\alpha}}) \left[ y_{i;\tilde{\alpha}} - \sum_j W_{ij} \phi_j(\mathbf{x}_{\tilde{\alpha}}) \right] \\ &= \sum_{\tilde{\alpha}} \phi_b(\mathbf{x}_{\tilde{\alpha}}) (z_{a;\tilde{\alpha}} - y_{a;\tilde{\alpha}}) = \sum_{\tilde{\alpha}} \varepsilon_{a;\tilde{\alpha}} \phi_b(\mathbf{x}_{\tilde{\alpha}}) \end{aligned} \quad (6)$$

where  $\varepsilon$  is the residual training error,

$$dW_{ij} = -\eta \sum_{\tilde{\alpha}} \phi_j(\mathbf{x}_{\tilde{\alpha}}) \varepsilon_{i;\tilde{\alpha}} . \quad (7)$$

So we have

$$dz_{i;\delta} = \sum_{a,b} \frac{\partial z_{i;\delta}}{\partial W_{ab}} dW_{ab} = -\eta \sum_{\tilde{\alpha}} k_{\delta \tilde{\alpha}} \varepsilon_{i;\tilde{\alpha}} . \quad (8)$$

In the linear model, the kernel  $k_{\delta \tilde{\alpha}}$  is static. The solution for the residual training error is

$$\varepsilon_{i;\tilde{\alpha}_1}(t) = \sum_{\tilde{\alpha}_2} U_{\tilde{\alpha}_1 \tilde{\alpha}_2}(t) \varepsilon_{i;\tilde{\alpha}_2}(0) , \quad (9)$$

where

$$\begin{aligned} U_{\tilde{\alpha}_t \tilde{\alpha}_0}(t) &\equiv [(1 - \eta k)^t]_{\tilde{\alpha}_t \tilde{\alpha}_0} \\ &= \sum_{\tilde{\alpha}_1, \dots, \tilde{\alpha}_{t-1}} (\delta_{\tilde{\alpha}_t \tilde{\alpha}_{t-1}} - \eta k_{\tilde{\alpha}_t \tilde{\alpha}_{t-1}}) \cdots (\delta_{\tilde{\alpha}_1 \tilde{\alpha}_0} - \eta k_{\tilde{\alpha}_1 \tilde{\alpha}_0}) . \end{aligned} \quad (10)$$

Now we could predict the model output for an arbitrary  $\delta$ . We have

$$\begin{aligned}
z_{i;\delta}(\infty) &= z_{i;\delta}(0) - \sum_{\tilde{\alpha}} k_{\delta\tilde{\alpha}} \left\{ \eta \sum_{s=0}^{\infty} \varepsilon_{i;\tilde{\alpha}}(s) \right\} \\
&= z_{i;\delta}(0) - \sum_{\tilde{\alpha}} k_{\delta\tilde{\alpha}} \left\{ \eta \sum_{s=0}^{\infty} \left[ \sum_{\tilde{\alpha}_1} U_{\tilde{\alpha}\tilde{\alpha}_1}(s) \varepsilon_{i;\tilde{\alpha}_1}(0) \right] \right\} \\
&= z_{i;\delta}(0) - \sum_{\tilde{\alpha}, \tilde{\alpha}_1} k_{\delta\tilde{\alpha}} \left\{ \eta \sum_{s=0}^{\infty} [(1 - \eta k)_{\tilde{\alpha}\tilde{\alpha}_1}^s] \varepsilon_{i;\tilde{\alpha}_1}(0) \right\} \\
&= z_{i;\delta}(0) - \sum_{\tilde{\alpha}, \tilde{\alpha}_1} k_{\delta\tilde{\alpha}} \left\{ \eta [1 - (1 - \eta k)]^{-1} \right\}_{\tilde{\alpha}\tilde{\alpha}_1} \varepsilon_{i;\tilde{\alpha}_1}(0) \\
&= z_{i;\delta}(0) - \sum_{\tilde{\alpha}, \tilde{\alpha}_1} k_{\delta\tilde{\alpha}} \tilde{k}^{\tilde{\alpha}\tilde{\alpha}_1} \varepsilon_{i;\tilde{\alpha}_1}(0), \tag{11}
\end{aligned}$$

where we have made use of geometric sums. It is easy to check that

$$z_{i;\tilde{\delta}}(\infty) = y_{i;\tilde{\delta}}, \tag{12}$$

where we take  $\tilde{\delta} \in \mathcal{A}$ .

## B. Quadratic model

Now we start to study perturbative corrections about the linear model. We consider the model definition,

$$z_{i;\delta}(\theta) = \sum_j W_{ij} \phi_j(\mathbf{x}_\delta) + \frac{\sigma}{2} \sum_{j_1, j_2} W_{ij_1} W_{ij_2} \psi_{j_1 j_2}(\mathbf{x}_\delta). \tag{13}$$

Here  $\sigma$  is a small number as a perturbative correction. Thus, the model prediction difference up to the quadratic order will be given by

$$dz_{i;\delta} = \sum_j \left[ \phi_j(\mathbf{x}_\delta) + \varepsilon \sum_{j_1=0} W_{ij_1} \psi_{j_1 j}(\mathbf{x}_\delta) \right] dW_{ij} + \frac{\sigma}{2} \sum_{j_1, j_2} \psi_{j_1 j_2}(\mathbf{x}_\delta) dW_{ij_1} dW_{ij_2}. \tag{14}$$

The first term here is the *effective feature map*,

$$\phi_{ij}^E(\mathbf{x}_\delta) \equiv \frac{dz_{i;\delta}}{dW_{ij}} = \phi_j(\mathbf{x}_\delta) + \sigma \sum_k W_{ik} \psi_{kj}(\mathbf{x}_\delta). \tag{15}$$

Now, we could collect our dynamical equations for  $z$ ,  $\phi^E$ ,  $W$  as

$$\begin{aligned}
dz_{i;\delta} &= \sum_j dW_{ij} \phi_{ij}^E(\mathbf{x}_\delta) + \frac{\sigma}{2} \sum_{j_1, j_2} dW_{ij_1} dW_{ij_2} \psi_{j_1 j_2}(\mathbf{x}_\delta), \\
d\phi_{ij}^E(\mathbf{x}_\delta) &= \sigma \sum_k dW_{ik} \psi_{kj}(\mathbf{x}_\delta), \\
dW_{ij} &= -\eta \sum_{\tilde{\alpha}} \phi_{ij}^E(\mathbf{x}_{\tilde{\alpha}}) \varepsilon_{i;\tilde{\alpha}}. \tag{16}
\end{aligned}$$

Moreover, now we have an MSE loss as

$$\mathcal{L}_{\mathcal{A}} = \frac{1}{2} \sum_{i, \tilde{\alpha}} \left[ y_{i;\tilde{\alpha}} - \sum_j W_{ij} \phi_j(\mathbf{x}_{\tilde{\alpha}}) - \frac{\sigma}{2} \sum_{j_1, j_2} W_{ij_1} W_{ij_2} \psi_{j_1 j_2}(\mathbf{x}_{\tilde{\alpha}}) \right]^2. \tag{17}$$

So we could write the dynamics of  $z$  more explicitly as

$$\begin{aligned} \dot{z}_{i;\delta} &= -\eta \sum_{\tilde{\alpha}} \left[ \sum_j \phi_{ij}^E(\mathbf{x}_\delta) \phi_{ij}^E(\mathbf{x}_{\tilde{\alpha}}) \right] \varepsilon_{i;\tilde{\alpha}} \\ &+ \frac{\eta^2}{2} \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} \left[ \sum_{j_1, j_2} \sigma \psi_{j_1 j_2}(\mathbf{x}_\delta) \phi_{ij_1}^E(\mathbf{x}_{\tilde{\alpha}_1}) \phi_{ij_2}^E(\mathbf{x}_{\tilde{\alpha}_2}) \right] \varepsilon_{i;\tilde{\alpha}_1} \varepsilon_{i;\tilde{\alpha}_2} . \end{aligned} \quad (18)$$

We define the *effective kernel*:

$$k_{ii;\delta_1 \delta_2}^E \equiv \sum_j \phi_{ij}^E(\mathbf{x}_{\delta_1}) \phi_{ij}^E(\mathbf{x}_{\delta_2}) , \quad (19)$$

and the *meta kernel*,

$$\begin{aligned} \mu_{\delta_0 \delta_1 \delta_2} &\equiv \sigma \sum_{j_1, j_2} \psi_{j_1 j_2}(\mathbf{x}_{\delta_0}) \phi_{j_1}(\mathbf{x}_{\delta_1}) \phi_{j_2}(\mathbf{x}_{\delta_2}) \\ &= \sigma \sum_{j_1, j_2} \varepsilon \psi_{j_1 j_2}(\mathbf{x}_{\delta_0}) \phi_{ij_1}^E(\mathbf{x}_{\delta_1}) \phi_{ij_2}^E(\mathbf{x}_{\delta_2}) + \mathcal{O}(\sigma^2) . \end{aligned} \quad (20)$$

So we have

$$\dot{z}_{i;\delta} = -\eta \sum_{\tilde{\alpha}} k_{ii;\delta \tilde{\alpha}}^E \varepsilon_{i;\tilde{\alpha}} + \frac{\eta^2}{2} \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} \mu_{\delta \tilde{\alpha}_1 \tilde{\alpha}_2} \varepsilon_{i;\tilde{\alpha}_1} \varepsilon_{i;\tilde{\alpha}_2} + \mathcal{O}(\sigma^2) . \quad (21)$$

The meta kernel is fixed. The effective kernel will satisfy the following dynamics

$$\dot{k}_{ii;\delta_1 \delta_2}^E = -\eta \sum_{\tilde{\alpha}} (\mu_{\delta_1 \delta_2 \tilde{\alpha}} + \mu_{\delta_2 \delta_1 \tilde{\alpha}}) \varepsilon_{i;\tilde{\alpha}} + \mathcal{O}(\sigma^2) . \quad (22)$$

Now we could try solving the output  $z$ . Based on perturbation theory, we could divide the whole output by the free term  $z^F$  and the interacting term  $z^I$ ,

$$z_{i;\delta}(t) \equiv z_{i;\delta}^F(t) + z_{i;\delta}^I(t) . \quad (23)$$

The free part follows the following exponential dynamics,

$$\dot{z}_{i;\delta}^F = -\eta \sum_{\tilde{\alpha}} k_{\delta \tilde{\alpha}} \varepsilon_{i;\tilde{\alpha}}^F \equiv z_{i;\delta}^F - \eta \sum_{\tilde{\alpha}} k_{\delta \tilde{\alpha}} [z_{i;\tilde{\alpha}}^F - y_{i;\tilde{\alpha}}] . \quad (24)$$

Here,  $k$  is the old definition of the kernel without quadratic terms, which is different from the effective kernel by  $O(\sigma)$ . Moreover, we have

$$\begin{aligned} k_{ii;\delta_1 \delta_2}^E(t) \\ = k_{ii;\delta_1 \delta_2}^E(0) - \sum_{\tilde{\alpha}} (\mu_{\delta_1 \delta_2 \tilde{\alpha}} + \mu_{\delta_2 \delta_1 \tilde{\alpha}}) a_{i;\tilde{\alpha}}(t) , \end{aligned} \quad (25)$$

where

$$a_{i;\tilde{\alpha}}(t) = \sum_{\tilde{\alpha}_2} \tilde{k}^{\tilde{\alpha} \tilde{\alpha}_2} (\varepsilon_{i;\tilde{\alpha}}(t) - \varepsilon_{i;\tilde{\alpha}}^F(t)) . \quad (26)$$

where  $\varepsilon^F$  is the free part of the residual training error  $z^F - y$ . Then, one could compute the interacting piece. We have

$$\dot{z}_{i;\delta}^I = - \sum_{j, \tilde{\alpha}} \eta k_{\delta \tilde{\alpha}} z_{j;\tilde{\alpha}}^I(t) + \eta \mathbb{F}_{i;\delta}(t) . \quad (27)$$

Here  $\mathbb{F}$  is the damping force

$$\mathbb{F}_{i;\delta}(t) \equiv - \sum_{\tilde{\alpha}} [k_{\delta\tilde{\alpha}}^E(t) - k_{\delta\tilde{\alpha}}] \varepsilon_{i;\tilde{\alpha}}^F(t) + \frac{\eta}{2} \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} \mu_{\delta\tilde{\alpha}_1\tilde{\alpha}_2} \varepsilon_{i;\tilde{\alpha}_1}^F(t) \varepsilon_{i;\tilde{\alpha}_2}^F(t), \quad (28)$$

and we have

$$z_{i;\delta}^I(t) = \eta \sum_{s=0}^{t-1} \left[ \mathbb{F}_{i;\delta}(s) - \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} k_{\delta\tilde{\alpha}_1} \tilde{k}^{\tilde{\alpha}_1\tilde{\alpha}_2} \mathbb{F}_{i;\tilde{\alpha}_2}(s) \right] + \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} k_{\delta\tilde{\alpha}_1} \tilde{k}^{\tilde{\alpha}_1\tilde{\alpha}_2} z_{i;\tilde{\alpha}_2}^I(t). \quad (29)$$

The whole system is a set of non-linear difference equations.

### 1. The residual training error

Now we could try to solve the residual training error. We note that, for  $\tilde{\alpha} \in \mathcal{A}$ , we have

$$\varepsilon_{i;\tilde{\alpha}}(t) - \varepsilon_{i;\tilde{\alpha}}^F(t) = z_{i;\tilde{\alpha}}^I(t). \quad (30)$$

So we have

$$k_{ii;\delta_1\delta_2}^E(t) = k_{ii;\delta_1\delta_2}^E(0) - \sum_{\tilde{\alpha}, \tilde{\alpha}_2} (\mu_{\delta_1\delta_2\tilde{\alpha}} + \mu_{\delta_2\delta_1\tilde{\alpha}}) \tilde{k}^{\tilde{\alpha}\tilde{\alpha}_2} z_{i;\tilde{\alpha}_2}^I(t). \quad (31)$$

Moreover, we call

$$k_{ii;\delta_1\delta_2}^\Delta \equiv k_{ii;\delta_1\delta_2}^\Delta \equiv k_{ii;\delta_1\delta_2}^E(0) - k_{\delta_1\delta_2}. \quad (32)$$

We note that the effective kernel looks like,

$$\begin{aligned} k_{ii;\delta_1\delta_2}^E &\equiv \sum_j \phi_{ij}^E(\mathbf{x}_{\delta_1}; 0) \phi_{ij}^E(\mathbf{x}_{\delta_2}; 0) \\ &= \sum_j \left( \phi_j(\mathbf{x}_{\delta_1}) + \sigma \sum_k W_{ik}(0) \psi_{kj}(\mathbf{x}_{\delta_1}) \right) \left( \phi_j(\mathbf{x}_{\delta_2}) + \sigma \sum_{k'} W_{ik'}(0) \psi_{k'j}(\mathbf{x}_{\delta_2}) \right) \\ &= k_{\delta_1\delta_2} + \sigma \sum_{j,k} W_{ik}(0) \psi_{kj}(\mathbf{x}_{\delta_1}) \phi_j(\mathbf{x}_{\delta_2}) + \sigma \sum_{j,k} W_{ik}(0) \psi_{kj}(\mathbf{x}_{\delta_2}) \phi_j(\mathbf{x}_{\delta_1}) + \mathcal{O}(\sigma^2). \end{aligned} \quad (33)$$

Keeping the leading order, we have

$$k_{ii;\delta_1\delta_2}^\Delta = \sigma \sum_{j,k} W_{ik}(0) \psi_{kj}(\mathbf{x}_{\delta_1}) \phi_j(\mathbf{x}_{\delta_2}) + \sigma \sum_{j,k} W_{ik}(0) \psi_{kj}(\mathbf{x}_{\delta_2}) \phi_j(\mathbf{x}_{\delta_1}), \quad (34)$$

and

$$\begin{aligned} k_{ii;\delta_1\delta_2}^E(t) &= k_{ii;\delta_1\delta_2}^E(0) - \sum_{\tilde{\alpha}, \tilde{\alpha}_2} (\mu_{\delta_1\delta_2\tilde{\alpha}} + \mu_{\delta_2\delta_1\tilde{\alpha}}) \tilde{k}^{\tilde{\alpha}\tilde{\alpha}_2} z_{i;\tilde{\alpha}_2}^I(t) \\ &= k_{\delta_1\delta_2} + k_{ii;\delta_1\delta_2}^\Delta - \sum_{\tilde{\alpha}, \tilde{\alpha}_2} (\mu_{\delta_1\delta_2\tilde{\alpha}} + \mu_{\delta_2\delta_1\tilde{\alpha}}) \tilde{k}^{\tilde{\alpha}\tilde{\alpha}_2} z_{i;\tilde{\alpha}_2}^I(t). \end{aligned} \quad (35)$$

Moreover, let us solve the damping force  $\mathbb{F}$ :

$$\begin{aligned}
\mathbb{F}_{i;\delta}(t) &= - \sum_{\tilde{\alpha}} [k_{\delta\tilde{\alpha}}^E(t) - k_{\delta\tilde{\alpha}}] \varepsilon_{i;\tilde{\alpha}}^F(t) + \frac{\eta}{2} \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} \mu_{\delta\tilde{\alpha}_1\tilde{\alpha}_2} \varepsilon_{i;\tilde{\alpha}_1}^F(t) \varepsilon_{i;\tilde{\alpha}_2}^F(t) \\
&= - \sum_{\tilde{\alpha}, \tilde{\alpha}_3} \left[ k_{ii;\delta\tilde{\alpha}}^E(0) - \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} (\mu_{\delta\tilde{\alpha}\tilde{\alpha}_1} + \mu_{\tilde{\alpha}\delta\tilde{\alpha}_1}) \tilde{k}^{\tilde{\alpha}_1\tilde{\alpha}_2} \varepsilon_{i;\tilde{\alpha}_2}^I(t) - k_{\delta\tilde{\alpha}} \right] U_{\tilde{\alpha}\tilde{\alpha}_3}(t) \varepsilon_{i;\tilde{\alpha}_3}(0) \\
&\quad + \frac{\eta}{2} \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4} \mu_{\delta\tilde{\alpha}_1\tilde{\alpha}_2} U_{\tilde{\alpha}_1\tilde{\alpha}_3}(t) U_{\tilde{\alpha}_2\tilde{\alpha}_4}(t) \varepsilon_{i;\tilde{\alpha}_3}(0) \varepsilon_{i;\tilde{\alpha}_4}(0) \\
&= - \sum_{\tilde{\alpha}, \tilde{\alpha}_3} k_{ii;\delta\tilde{\alpha}}^\Delta U_{\tilde{\alpha}\tilde{\alpha}_3}(t) \varepsilon_{i;\tilde{\alpha}_3}(0) + \sum_{\tilde{\alpha}, \tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3} (\mu_{\delta\tilde{\alpha}\tilde{\alpha}_1} + \mu_{\tilde{\alpha}\delta\tilde{\alpha}_1}) \tilde{k}^{\tilde{\alpha}_1\tilde{\alpha}_2} U_{\tilde{\alpha}\tilde{\alpha}_3}(t) \varepsilon_{i;\tilde{\alpha}_3}(0) z_{i;\tilde{\alpha}_2}^I(t) \\
&\quad + \frac{\eta}{2} \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3, \tilde{\alpha}_4} \mu_{\delta\tilde{\alpha}_1\tilde{\alpha}_2} U_{\tilde{\alpha}_1\tilde{\alpha}_3}(t) U_{\tilde{\alpha}_2\tilde{\alpha}_4}(t) \varepsilon_{i;\tilde{\alpha}_3}(0) \varepsilon_{i;\tilde{\alpha}_4}(0) .
\end{aligned} \tag{36}$$

The second and the last term in the last equality contributes higher orders. Thus, if the initial  $k^\Delta$  is not vanishing, we have

$$\mathbb{F}_{i;\delta}(t) = - \sum_{\tilde{\alpha}, \tilde{\alpha}_3} k_{ii;\delta\tilde{\alpha}}^\Delta U_{\tilde{\alpha}\tilde{\alpha}_3}(t) \varepsilon_{i;\tilde{\alpha}_3}(0) + \mathcal{O}(\eta\sigma) . \tag{37}$$

The contribution in the first term will dominate at least in the early time when  $k^\Delta$  is not vanishing. In the late time where  $U$  decays significantly, we would have some non-perturbative effects.

Thus, we can plug the expression back to solve  $z^I(t)$ . We have

$$\begin{aligned}
z_{i;\tilde{\alpha}}^I(t) &= \eta \sum_{s=0}^{t-1} \sum_{\tilde{\alpha}_1} U_{\tilde{\alpha}\tilde{\alpha}_1}(t-1-s) \mathbb{F}_{i;\tilde{\alpha}_1}(s) \\
&= -\eta \sum_{s=0}^{t-1} \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2, \tilde{\alpha}_3} U_{\tilde{\alpha}\tilde{\alpha}_1}(t-1-s) U_{\tilde{\alpha}_2\tilde{\alpha}_3}(s) k_{ii;\tilde{\alpha}_1\tilde{\alpha}_2}^\Delta \varepsilon_{i;\tilde{\alpha}_3}(0) \\
&= -\eta \sum_{s=0}^{t-1} (1-\eta k)^{t-1-s} k_{ii}^\Delta (\delta - \eta k)^s \varepsilon_i(0) ,
\end{aligned} \tag{38}$$

The last formula is given in the following matrix form:

$$\begin{aligned}
(1-\eta k)_{\tilde{\alpha}_1\tilde{\alpha}_2} &= \delta_{\tilde{\alpha}_1\tilde{\alpha}_2} - \eta k_{\tilde{\alpha}_1\tilde{\alpha}_2} = \delta_{\tilde{\alpha}_1\tilde{\alpha}_2} - \eta \sum_j \phi_j(\mathbf{x}_{\tilde{\alpha}_1}) \phi_j(\mathbf{x}_{\tilde{\alpha}_2}) , \\
(k_{ii}^\Delta)_{\tilde{\alpha}_1\tilde{\alpha}_2} &= \sigma \sum_{j,k} W_{ik}(0) \psi_{kj}(\mathbf{x}_{\tilde{\alpha}_1}) \phi_j(\mathbf{x}_{\tilde{\alpha}_2}) + \sigma \sum_{j,k} W_{ik}(0) \psi_{kj}(\mathbf{x}_{\tilde{\alpha}_2}) \phi_j(\mathbf{x}_{\tilde{\alpha}_1}) \equiv \sigma (\mathbf{M}^i)_{\tilde{\alpha}_1, \tilde{\alpha}_2} , \\
(\varepsilon_i(0))_{\tilde{\alpha}} &= \varepsilon_{i;\tilde{\alpha}}(0) .
\end{aligned} \tag{39}$$

Moreover, we could indeed get more information by just making the bound. We have

$$\begin{aligned}
\|z_{i;\tilde{\alpha}}^I(t)\| &= \left\| \eta \sum_{s=0}^{t-1} (1-\eta k)^{t-1-s} k_{ii}^\Delta (1-\eta k)^s \varepsilon_i(0) \right\| \\
&\leq \eta \sum_{s=0}^{t-1} \|1-\eta k\|^{t-1-s} \|k_{ii}^\Delta\| \|1-\eta k\|^s \|\varepsilon_i(0)\| \\
&= \eta \left( \sum_{s=0}^{t-1} \|1-\eta k\|^{t-1-s} \right) \|k_{ii}^\Delta\| \|\varepsilon_i(0)\| \\
&= \eta t \|1-\eta k\|^{t-1} \|k_{ii}^\Delta\| \|\varepsilon_i(0)\| = \sigma \eta t \|1-\eta k\|^{t-1} \|\mathbf{M}^i\| \|\varepsilon_i(0)\| .
\end{aligned} \tag{40}$$

Now we compare the convergence time noticing that

$$\varepsilon_{i;\tilde{\alpha}_1}^F(t) = \sum_{\tilde{\alpha}_2} U_{\tilde{\alpha}_1\tilde{\alpha}_2}(t) \varepsilon_{i;\tilde{\alpha}_2}(0) . \tag{41}$$

So

$$\|\varepsilon_i^F(t)\| \leq \|1 - \eta k\|^t \|\varepsilon_i(0)\|. \quad (42)$$

Schematically, we have

$$\frac{\|z_i^I(t)\|}{\|\varepsilon_i^F(t)\|} \sim \frac{\sigma \eta t \|\mathbf{M}^i\| \|\varepsilon_i(0)\|}{\|1 - \eta k\| \|\varepsilon_i(0)\|} \sim \frac{\sigma \eta t \|\mathbf{M}^i\|}{\|1 - \eta k\|}. \quad (43)$$

The relative perturbative error contribution will grow linearly in time.

## 2. The asymptotic regime

Now instead of only looking at the training set, we study the asymptotic regime for general inputs. We start from

$$z_{i;\delta}^I(t) = \eta \sum_{s=0}^{t-1} \left[ \mathbb{F}_{i;\delta}(s) - \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} k_{\delta \tilde{\alpha}_1} \tilde{k}^{\tilde{\alpha}_1 \tilde{\alpha}_2} \mathbb{F}_{i;\tilde{\alpha}_2}(s) \right] + \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} k_{\delta \tilde{\alpha}_1} \tilde{k}^{\tilde{\alpha}_1 \tilde{\alpha}_2} z_{i;\tilde{\alpha}_2}^I(t). \quad (44)$$

At the asymptotic convergence, the interacting perturbative correction on the training set will converge to zero, so we have

$$z_{i;\delta}^I(\infty) = \left[ \eta \sum_{s=0}^{\infty} \mathbb{F}_{i;\delta}(s) \right] - \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} k_{\delta \tilde{\alpha}_1} \tilde{k}^{\tilde{\alpha}_1 \tilde{\alpha}_2} \left[ \eta \sum_{s=0}^{\infty} \mathbb{F}_{i;\tilde{\alpha}_2}(s) \right]. \quad (45)$$

So we need to perform the sum

$$\begin{aligned} \eta \sum_{s=0}^{\infty} \mathbb{F}_{i;\delta}(s) &= \sum_{\tilde{\alpha}_0, \tilde{\alpha}} (\mu_{\delta \tilde{\alpha} \tilde{\alpha}_0} + \mu_{\tilde{\alpha} \delta \tilde{\alpha}_0}) \left\{ \eta \sum_{s=0}^{\infty} a_{i;\tilde{\alpha}_0}(s) \varepsilon_{i;\tilde{\alpha}}^F(s) \right\} \\ &+ \frac{\eta}{2} \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} \mu_{\delta \tilde{\alpha}_1 \tilde{\alpha}_2} \left\{ \eta \sum_{s=0}^{\infty} \varepsilon_{i;\tilde{\alpha}_1}^F(s) \varepsilon_{i;\tilde{\alpha}_2}^F(s) \right\}. \end{aligned} \quad (46)$$

We note that

$$\begin{aligned} \eta \sum_{s=0}^{\infty} \varepsilon_{i;\tilde{\alpha}_1}^F(t) \varepsilon_{i;\tilde{\alpha}_2}^F(t) &= \eta \sum_{\tilde{\alpha}_3, \tilde{\alpha}_4} \sum_{s=0}^{\infty} [(1 - \eta k)^s]_{\tilde{\alpha}_1 \tilde{\alpha}_3} [(1 - \eta k)^s]_{\tilde{\alpha}_2 \tilde{\alpha}_4} \varepsilon_{i;\tilde{\alpha}_3}(0) \varepsilon_{i;\tilde{\alpha}_4}(0) \\ &= \sum_{\tilde{\alpha}_3, \tilde{\alpha}_4} X_{\parallel}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \varepsilon_{i;\tilde{\alpha}_3}(0) \varepsilon_{i;\tilde{\alpha}_4}(0). \end{aligned} \quad (47)$$

Here, we define the following inverting tensor:

$$X_{\parallel}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} = \sum_{s=0}^{\infty} [(1 - \eta k)^s]_{\tilde{\alpha}_1 \tilde{\alpha}_3} [(1 - \eta k)^s]_{\tilde{\alpha}_2 \tilde{\alpha}_4}, \quad (48)$$

which is implicitly defined as

$$\begin{aligned} \delta_{\tilde{\alpha}_5}^{\tilde{\alpha}_1} \delta_{\tilde{\alpha}_6}^{\tilde{\alpha}_2} &= \sum_{\tilde{\alpha}_3, \tilde{\alpha}_4} X_{\parallel}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \frac{1}{\eta} \left[ \delta_{\tilde{\alpha}_3 \tilde{\alpha}_5} \delta_{\tilde{\alpha}_4 \tilde{\alpha}_6} - (\delta_{\tilde{\alpha}_3 \tilde{\alpha}_5} - \eta \tilde{k}_{\tilde{\alpha}_3 \tilde{\alpha}_5}) (\delta_{\tilde{\alpha}_4 \tilde{\alpha}_6} - \eta \tilde{k}_{\tilde{\alpha}_4 \tilde{\alpha}_6}) \right] \\ &= \sum_{\tilde{\alpha}_3, \tilde{\alpha}_4} X_{\parallel}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \left( \tilde{k}_{\tilde{\alpha}_3 \tilde{\alpha}_5} \delta_{\tilde{\alpha}_4 \tilde{\alpha}_6} + \delta_{\tilde{\alpha}_3 \tilde{\alpha}_5} \tilde{k}_{\tilde{\alpha}_4 \tilde{\alpha}_6} - \eta \tilde{k}_{\tilde{\alpha}_3 \tilde{\alpha}_5} \tilde{k}_{\tilde{\alpha}_4 \tilde{\alpha}_6} \right). \end{aligned} \quad (49)$$

Using the inverting tensor, the final expression is given by

$$\begin{aligned} z_{i;\delta}(\infty) &= z_{i;\delta}(0) - \sum_{\tilde{\alpha}_1, \tilde{\alpha}_2} k_{\delta \tilde{\alpha}_1} \tilde{k}^{\tilde{\alpha}_1 \tilde{\alpha}_2} \varepsilon_{i;\tilde{\alpha}_2}(0) \\ &+ \sum_{\tilde{\alpha}_1, \dots, \tilde{\alpha}_4} \left[ \mu_{\tilde{\alpha}_1 \delta \tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6 \in \mathcal{A}} k_{\delta \tilde{\alpha}_5} \tilde{k}^{\tilde{\alpha}_5 \tilde{\alpha}_6} \mu_{\tilde{\alpha}_1 \tilde{\alpha}_6 \tilde{\alpha}_2} \right] Z_A^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \varepsilon_{i;\tilde{\alpha}_3}(0) \varepsilon_{i;\tilde{\alpha}_4}(0) \\ &+ \sum_{\tilde{\alpha}_1, \dots, \tilde{\alpha}_4} \left[ \mu_{\delta \tilde{\alpha}_1 \tilde{\alpha}_2} - \sum_{\tilde{\alpha}_5, \tilde{\alpha}_6 \in \mathcal{A}} k_{\delta \tilde{\alpha}_5} \tilde{k}^{\tilde{\alpha}_5 \tilde{\alpha}_6} \mu_{\tilde{\alpha}_6 \tilde{\alpha}_1 \tilde{\alpha}_2} \right] Z_B^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \varepsilon_{i;\tilde{\alpha}_3}(0) \varepsilon_{i;\tilde{\alpha}_4}(0), \end{aligned} \quad (50)$$

where  $Z$ s are the *algorithm projectors*,

$$\begin{aligned} Z_A^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} &\equiv \tilde{k}^{\tilde{\alpha}_1 \tilde{\alpha}_3} \tilde{k}^{\tilde{\alpha}_2 \tilde{\alpha}_4} - \sum_{\tilde{\alpha}_5} \tilde{k}^{\tilde{\alpha}_2 \tilde{\alpha}_5} X_{\parallel}^{\tilde{\alpha}_1 \tilde{\alpha}_5 \tilde{\alpha}_3 \tilde{\alpha}_4}, \\ Z_B^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} &\equiv \tilde{k}^{\tilde{\alpha}_1 \tilde{\alpha}_3} \tilde{k}^{\tilde{\alpha}_2 \tilde{\alpha}_4} - \sum_{\tilde{\alpha}_5} \tilde{k}^{\tilde{\alpha}_2 \tilde{\alpha}_5} X_{\parallel}^{\tilde{\alpha}_1 \tilde{\alpha}_5 \tilde{\alpha}_3 \tilde{\alpha}_4} + \frac{\eta}{2} X_{\parallel}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4}. \end{aligned} \quad (51)$$

Finally, in the continuum limit, we could drop out the  $\eta$  terms in the algorithm projectors,

$$\begin{aligned} Z_A^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} &= Z_B^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} \equiv \tilde{k}^{\tilde{\alpha}_1 \tilde{\alpha}_3} \tilde{k}^{\tilde{\alpha}_2 \tilde{\alpha}_4} - \sum_{\tilde{\alpha}_5} \tilde{k}^{\tilde{\alpha}_2 \tilde{\alpha}_5} X_{\parallel}^{\tilde{\alpha}_1 \tilde{\alpha}_5 \tilde{\alpha}_3 \tilde{\alpha}_4}, \\ \sum_{\tilde{\alpha}_3, \tilde{\alpha}_4} X_{\parallel}^{\tilde{\alpha}_1 \tilde{\alpha}_2 \tilde{\alpha}_3 \tilde{\alpha}_4} &\left( \tilde{k}_{\tilde{\alpha}_3 \tilde{\alpha}_5} \delta_{\tilde{\alpha}_4 \tilde{\alpha}_6} + \delta_{\tilde{\alpha}_3 \tilde{\alpha}_5} \tilde{k}_{\tilde{\alpha}_4 \tilde{\alpha}_6} \right) = \delta_{\tilde{\alpha}_5}^{\tilde{\alpha}_1} \delta_{\tilde{\alpha}_6}^{\tilde{\alpha}_2}. \end{aligned} \quad (52)$$

The existence of algorithm projectors shows algorithm dependence in those perturbative corrections. This will typically happen when the model has multiple (local) minimal.

## II. DETAILS ON THE QUANTUM OPTIMIZATION

### A. General setup

The quantum optimization problem we discuss here has a simpler formulation compared to quantum machine learning. Since the loss function does not contain the training data, we do not need to consider the difference between the training set and the whole space. It could be understood as a limiting case of the quantum machine learning problem.

We use the loss function,

$$\mathcal{L}(\theta) = \frac{1}{2} (\langle \Psi_0 | U^\dagger O U | \Psi_0 \rangle - O_0)^2 \equiv \frac{1}{2} \varepsilon^2. \quad (53)$$

The trainable ansatz is,

$$U = \prod_{\ell=1}^L W_\ell U_\ell = \prod_{\ell=1}^L W_\ell \exp(i\theta_\ell X_\ell). \quad (54)$$

Here, note that since the operator  $O$  is Hermitian, the loss function is always real and non-negative. The gradient descent is

$$d\theta_\ell = -\eta \frac{d\mathcal{L}(\theta)}{d\theta_\ell} = -\eta (\langle \Psi_0 | U^\dagger O U | \Psi_0 \rangle - O_0) \frac{d\langle \Psi_0 | U^\dagger O U | \Psi_0 \rangle}{d\theta_\ell} = -\eta \varepsilon \frac{d\varepsilon}{d\theta_\ell}, \quad (55)$$

where

$$\varepsilon = \langle \Psi_0 | U^\dagger O U | \Psi_0 \rangle - O_0. \quad (56)$$

Thus

$$d\varepsilon = \sum_\ell \frac{d\varepsilon}{d\theta_\ell} d\theta_\ell = -\eta \sum_\ell \frac{d\varepsilon}{d\theta_\ell} \frac{d\varepsilon}{d\theta_\ell} \varepsilon. \quad (57)$$

The object

$$K = \sum_\ell \frac{d\varepsilon}{d\theta_\ell} \frac{d\varepsilon}{d\theta_\ell}, \quad (58)$$

is the optimization analog of the QNTK. More precisely, we have

$$\frac{d\varepsilon}{d\theta_\ell} = \frac{d\langle \Psi_0 | U^\dagger O U | \Psi_0 \rangle}{d\theta_\ell} = -i \langle \Psi_0 | U_{+, \ell}^\dagger [X_\ell, U_\ell^\dagger W_\ell^\dagger U_{-, \ell}^\dagger O U_{-, \ell} W_\ell U_\ell] U_{+, \ell} | \Psi_0 \rangle, \quad (59)$$

with the help of the definition,

$$\begin{aligned} U_{-, \ell} &\equiv \prod_{\ell'=1}^{\ell-1} W_{\ell'} U_{\ell'} , \\ U_{+, \ell} &\equiv \prod_{\ell'=\ell+1}^L W_{\ell'} U_{\ell'} . \end{aligned} \quad (60)$$

So,

$$d\varepsilon = -\eta \sum_{\ell} \frac{d\varepsilon}{d\theta_{\ell}} \frac{d\varepsilon}{d\theta_{\ell}} \varepsilon = \eta \varepsilon \left\langle \Psi_0 \left| U_{+, \ell}^{\dagger} \left[ X_{\ell}, U_{\ell}^{\dagger} W_{\ell}^{\dagger} U_{-, \ell}^{\dagger} O U_{-, \ell} W_{\ell} U_{\ell} \right] U_{+, \ell} \right| \Psi_0 \right\rangle^2 . \quad (61)$$

Can we solve the difference equation? We have

$$d\theta_{\ell} = i\eta \left\langle \Psi_0 \left| U_{+, \ell}^{\dagger} \left[ X_{\ell}, U_{\ell}^{\dagger} W_{\ell}^{\dagger} U_{-, \ell}^{\dagger} O U_{-, \ell} W_{\ell} U_{\ell} \right] U_{+, \ell} \right| \Psi_0 \right\rangle . \quad (62)$$

In the continuum limit, we have

$$\frac{d\theta_{\ell}}{dt} = i\eta_c \left\langle \Psi_0 \left| U_{+, \ell}^{\dagger} \left[ X_{\ell}, U_{\ell}^{\dagger} W_{\ell}^{\dagger} U_{-, \ell}^{\dagger} O U_{-, \ell} W_{\ell} U_{\ell} \right] U_{+, \ell} \right| \Psi_0 \right\rangle , \quad (63)$$

where  $\eta_c = \eta/dt$  is the continuous version of the learning rate. In principle, this is a coupled nonlinear ODE system, and one could use the ODE theory to study them. Moreover, we have

$$U_{+, \ell}^{\dagger} \left[ X_{\ell}, U_{\ell}^{\dagger} W_{\ell}^{\dagger} U_{-, \ell}^{\dagger} O U_{-, \ell} W_{\ell} U_{\ell} \right] U_{+, \ell} = U_{+, \ell}^{\dagger} U_{\ell}^{\dagger} \left[ X_{\ell}, W_{\ell}^{\dagger} U_{-, \ell}^{\dagger} O U_{-, \ell} W_{\ell} \right] U_{\ell} U_{+, \ell} . \quad (64)$$

Defining

$$A_{\ell} = \left[ X_{\ell}, W_{\ell}^{\dagger} U_{-, \ell}^{\dagger} O U_{-, \ell} W_{\ell} \right] , \quad (65)$$

and if we are expanding  $\theta_{\ell}$  when it is slightly deviated from 0 (more generally, around a fixed angle, which is equivalent to a redefinition of constant gates), we have,

$$U_{+, \ell}^{\dagger} U_{\ell}^{\dagger} A_{\ell} U_{\ell} U_{+, \ell} = U_{+, \ell}^{\dagger} A_{\ell} U_{+, \ell} - i\theta_{\ell} U_{+, \ell}^{\dagger} [X_{\ell}, A_{\ell}] U_{+, \ell} + \mathcal{O}(\theta_{\ell}^2) . \quad (66)$$

Thus, we note that the second-order expansions of  $\theta_{\ell}$  here means second-order commutators among  $X_{\ell}$  and some unitary-addressed versions of the operator  $O$ , which is the operator we wish to optimize. Our experiences could be easily generalized to high orders.

## B. Frozen QNTK

Now, we rescale the original variational angles by a factor  $\delta$  where the variational angles are around  $\theta^* + \delta\varphi$ . The constant term  $\theta^*$  will produce a constant gate  $\exp(i\theta_{\ell}^* X_{\ell})$ , such that  $W_{\ell} U_{\ell} \rightarrow W_{\ell} \exp(i\theta_{\ell}^* X_{\ell}) \exp(i\delta\varphi_{\ell} X_{\ell})$ . So we could absorb the  $\theta^*$  dependence to the definition of the constant gate by defining  $W_{\ell}(\theta_{\ell}^*) \equiv W_{\ell} \exp(i\theta_{\ell}^* X_{\ell})$ , and we have  $W_{\ell} U_{\ell} \rightarrow W_{\ell}(\theta_{\ell}^*) \exp(i\delta\varphi_{\ell} X_{\ell})$ . Thus, for simplicity, we could understand the trainable gate as  $U_{\ell} \rightarrow \exp(i\delta\varphi_{\ell} X_{\ell})$  with a redefined constant gate. In our calculation later, we will drop the  $\theta^*$  notation and understand the variational angles as small parameters rescaled by  $\delta$  for our notation convenience.

In the frozen QNTK limit, we have,

$$K = -\delta^2 \sum_{\ell} \left\langle \Psi_0 \left| W_{+, \ell}^{\dagger} \left[ X_{\ell}, W_{\ell}^{\dagger} W_{-, \ell}^{\dagger} O W_{-, \ell} W_{\ell} \right] W_{+, \ell} \right| \Psi_0 \right\rangle^2 . \quad (67)$$

And we define

$$W_{-, \ell} \equiv \prod_{\ell'=1}^{\ell-1} W_{\ell'} , \quad W_{+, \ell} \equiv \prod_{\ell'=\ell+1}^L W_{\ell'} . \quad (68)$$

Note that the frozen QNTK does not depend on the variational parameters. The gradient descent dynamics, in this case, is very easy to solve. We have

$$\varepsilon(t) = (1 - \eta K)^t \varepsilon(0) \equiv \left( 1 + \eta \delta^2 \sum_{\ell} \left\langle \Psi_0 \left| W_{+, \ell}^{\dagger} \left[ X_{\ell}, W_{\ell}^{\dagger} W_{-, \ell}^{\dagger} O W_{-, \ell} W_{\ell} \right] W_{+, \ell} \right| \Psi_0 \right\rangle^2 \right)^t \varepsilon(0). \quad (69)$$

The convergence rate is given by,

$$\begin{aligned} \tau_c &= -\log(1 - \eta K) \approx \eta K \\ &= \eta \delta^2 \sum_{\ell} \left\langle \Psi_0 \left| W_{+, \ell}^{\dagger} \left[ X_{\ell}, W_{\ell}^{\dagger} W_{-, \ell}^{\dagger} O W_{-, \ell} W_{\ell} \right] W_{+, \ell} \right| \Psi_0 \right\rangle^2 \\ &\leq 2\eta \delta^2 L \|O\|_{\max}^2 \|X_{\ell}\|^2. \end{aligned} \quad (70)$$

### C. dQNTK

Now let us focus on the second order to develop an analog of representation learning theory in the quantum optimization example. We have

$$\begin{aligned} \varepsilon &= \left\langle \Psi_0 \left| \left( \prod_{\ell'=L}^1 W_{\ell'}^{\dagger} \right) O \left( \prod_{\ell=1}^L W_{\ell} \right) \right| \Psi_0 \right\rangle - O_0 \\ &\quad - i\delta \sum_{\ell} \varphi_{\ell} \left\langle \Psi_0 \left| W_{+, \ell}^{\dagger} \left[ X_{\ell}, W_{\ell}^{\dagger} W_{-, \ell}^{\dagger} O W_{-, \ell} W_{\ell} \right] W_{+, \ell} \right| \Psi_0 \right\rangle \\ &\quad - \frac{\delta^2}{2} \sum_{\ell_1, \ell_2} \varphi_{\ell_1} \varphi_{\ell_2} \times \begin{cases} \left\langle \Psi_0 \left| W_{+, \ell_1}^{\dagger} \left[ X_{\ell_1}, Q_{\ell_1, \ell_2}^{\dagger} \left[ X_{\ell_2}, W_{\ell_2}^{\dagger} W_{-, \ell_2}^{\dagger} O W_{-, \ell_2} W_{\ell_2} \right] Q_{\ell_2, \ell_1} \right] W_{+, \ell_1} \right| \Psi_0 \right\rangle : \ell_1 \geq \ell_2 \\ \left\langle \Psi_0 \left| W_{+, \ell_2}^{\dagger} \left[ X_{\ell_2}, Q_{\ell_2, \ell_1}^{\dagger} \left[ X_{\ell_1}, W_{\ell_1}^{\dagger} W_{-, \ell_1}^{\dagger} O W_{-, \ell_1} W_{\ell_1} \right] Q_{\ell_1, \ell_2} \right] W_{+, \ell_2} \right| \Psi_0 \right\rangle : \ell_1 < \ell_2 \end{cases}, \end{aligned} \quad (71)$$

where

$$\begin{aligned} W_{\ell_1, \ell_2} &\equiv \prod_{\ell=\ell_1+1}^{\ell_2-1} W_{\ell}, \\ Q_{\ell_1, \ell_2} &= \begin{cases} W_{\ell_1, \ell_2} W_{\ell_2} : \ell_1 < \ell_2 \\ 1 : \ell_1 = \ell_2 \end{cases}. \end{aligned} \quad (72)$$

Moreover, we have

$$d\varphi_{\ell} = -\eta \varepsilon \frac{d\varepsilon}{d\varphi_{\ell}}. \quad (73)$$

Thus

$$d\varepsilon = -\eta \sum_{\ell} \frac{d\varepsilon}{d\varphi_{\ell}} \frac{d\varepsilon}{d\varphi_{\ell}} \varepsilon + \frac{1}{2} \eta^2 \varepsilon^2 \sum_{\ell_1, \ell_2} \frac{d^2\varepsilon}{d\varphi_{\ell_1} d\varphi_{\ell_2}} \frac{d\varepsilon}{d\varphi_{\ell_1}} \frac{d\varepsilon}{d\varphi_{\ell_2}}. \quad (74)$$

The structure of the gradient descent equation for the residual optimization error is very similar to the case in the quadratic model and representation learning context. We define the free part and the interacting part of  $\varepsilon$  as

$$\varepsilon = \varepsilon^F + \varepsilon^I. \quad (75)$$

The free part is given by,

$$\begin{aligned} \varepsilon^F &= (1 - \eta K)^t \varepsilon(0), \\ K &= -\delta^2 \sum_{\ell} \left\langle \Psi_0 \left| W_{+, \ell}^{\dagger} \left[ X_{\ell}, W_{\ell}^{\dagger} W_{-, \ell}^{\dagger} O W_{-, \ell} W_{\ell} \right] W_{+, \ell} \right| \Psi_0 \right\rangle^2, \end{aligned} \quad (76)$$

and the interacting part is given by

$$\varepsilon^I(t) = -\eta t(1 - \eta K)^{t-1} K^\Delta \varepsilon(0), \quad (77)$$

where we define the effective kernel up to the dQNTK order,

$$K^E = \sum_{\ell} \frac{d\varepsilon}{d\varphi_{\ell}} \frac{d\varepsilon}{d\varphi_{\ell}}. \quad (78)$$

And we have

$$\begin{aligned} K^\Delta &\equiv K^E(0) - K = \left( \sum_{\ell} \frac{d\varepsilon}{d\theta_{\ell}} \frac{d\varepsilon}{d\theta_{\ell}} \right) (0) - \sum_{\ell} \frac{d\varepsilon^F}{d\theta_{\ell}} \frac{d\varepsilon^F}{d\theta_{\ell}} \\ &= 2i\delta^3 \sum_{\ell} \left\langle \Psi_0 \left| W_{+, \ell}^\dagger \left[ X_{\ell}, W_{\ell}^\dagger W_{-, \ell}^\dagger O W_{-, \ell} W_{\ell} \right] W_{+, \ell} \right| \Psi_0 \right\rangle \\ &\quad \sum_{\ell'} \left\langle \Psi_0 \left| W_{+, \ell'}^\dagger \left[ X_{\ell'}, Q_{\ell', \ell}^\dagger \left[ X_{\ell}, W_{\ell}^\dagger W_{-, \ell}^\dagger O W_{-, \ell} W_{\ell} \right] W_{\ell, \ell'} W_{\ell'} \right] W_{+, \ell'} \right| \Psi_0 \right\rangle : \ell' \geq \ell \\ &\quad \left\langle \Psi_0 \left| W_{+, \ell}^\dagger \left[ X_{\ell}, Q_{\ell, \ell'}^\dagger \left[ X_{\ell'}, W_{\ell'}^\dagger W_{-, \ell'}^\dagger O W_{-, \ell'} W_{\ell'} \right] W_{\ell', \ell} W_{\ell} \right] W_{+, \ell} \right| \Psi_0 \right\rangle : \ell' < \ell \quad \varphi_{\ell'}(0). \end{aligned} \quad (79)$$

Similarly, one define the quantum meta-kernel (dQNTK) as

$$\begin{aligned} \mu &= \delta^4 \sum_{\ell_1, \ell_2} \left\langle \Psi_0 \left| W_{+, \ell_1}^\dagger \left[ X_{\ell_1}, W_{\ell_1}^\dagger W_{-, \ell_1}^\dagger O W_{-, \ell_1} W_{\ell_1} \right] W_{+, \ell_1} \right| \Psi_0 \right\rangle \\ &\quad \left\langle \Psi_0 \left| W_{+, \ell_2}^\dagger \left[ X_{\ell_2}, W_{\ell_2}^\dagger W_{-, \ell_2}^\dagger O W_{-, \ell_2} W_{\ell_2} \right] W_{+, \ell_2} \right| \Psi_0 \right\rangle \times \\ &\quad \left( \left\langle \Psi_0 \left| W_{+, \ell_1}^\dagger \left[ X_{\ell_1}, Q_{\ell_1, \ell_2}^\dagger \left[ X_{\ell_2}, W_{\ell_2}^\dagger W_{-, \ell_2}^\dagger O W_{-, \ell_2} W_{\ell_2} \right] Q_{\ell_2, \ell_1} \right] W_{+, \ell_1} \right| \Psi_0 \right\rangle : \ell_1 \geq \ell_2 \right) \\ &\quad \left( \left\langle \Psi_0 \left| W_{+, \ell_2}^\dagger \left[ X_{\ell_2}, Q_{\ell_2, \ell_1}^\dagger \left[ X_{\ell_1}, W_{\ell_1}^\dagger W_{-, \ell_1}^\dagger O W_{-, \ell_1} W_{\ell_1} \right] Q_{\ell_1, \ell_2} \right] W_{+, \ell_2} \right| \Psi_0 \right\rangle : \ell_1 < \ell_2 \right). \end{aligned} \quad (80)$$

Finally, we wish to mention that in [1], for the classical neural networks they study, the leading order perturbative contribution  $\mathcal{O}(1/\text{width})$  is both given by dNTK and ddNTK in dynamics. In our frozen QNTK limit (in the context of lazy training), this does not happen because of power counting in  $\delta$ .

Finally, we wish to mention that the scaling of  $\varepsilon^I$  will lead to the so-called *catapult effect*. In general, in higher order corrections, we get schematically the correction  $\sim t^p \exp(-\eta K t)$  for a more general polynomial  $t^p$  in the prefactor of the exponential decay. This type of correction forms a first-principle explanation of the catapult effect, where a similar related model is discussed. A full characterization of the catapult effect in classical and quantum cases is beyond the scope of this paper, and we refer the detailed research to the future [4].

### III. DETAILS ON THE QUANTUM MACHINE LEARNING: HERMITIAN OPERATOR EXPECTATION VALUE EVALUATION

#### A. General setup

We define our model as

$$z_{i;\delta} \equiv z_i(\theta, \mathbf{x}_\delta) = \langle \phi(\mathbf{x}_\delta) | U^\dagger O_i U | \phi(\mathbf{x}_\delta) \rangle. \quad (81)$$

where  $O_i$  is the  $i$ -th Hermitian observable. We assume that  $O_i$  is taken from a subset of Hermitian operators of the Hilbert space  $\mathcal{H}$ , denoted by  $\mathcal{O}(\mathcal{H})$ . The dimension of  $\mathcal{O}(\mathcal{H})$  is upper bounded by polynomials of the dimension of the Hilbert space,  $\dim \mathcal{H}$ . The trainable ansatz is,

$$U = \prod_{\ell=1}^L W_{\ell} U_{\ell} = \prod_{\ell=1}^L W_{\ell} \exp(i\theta_{\ell} X_{\ell}). \quad (82)$$

One could take the derivative,

$$dz_{i;\delta} = \sum_{\ell} \frac{dz_{i;\delta}}{d\theta_{\ell}} d\theta_{\ell}, \quad (83)$$

The loss is

$$L_{\mathcal{A}}(\theta) = \frac{1}{2} \sum_{\bar{\alpha}, i} (y_{i; \bar{\alpha}} - z_{i; \bar{\alpha}})^2 = \frac{1}{2} \sum_{\bar{\alpha}, i} \varepsilon_{i; \bar{\alpha}}^2. \quad (84)$$

So

$$\frac{dL_{\mathcal{A}}(\theta)}{d\theta_{\ell}} = \sum_{\bar{\alpha}, i} \varepsilon_{i; \bar{\alpha}} \frac{dz_{i; \bar{\alpha}}}{d\theta_{\ell}}. \quad (85)$$

The gradient descent rule is

$$d\theta_{\ell} = -\eta \frac{dL_{\mathcal{A}}}{d\theta_{\ell}} = -\eta \sum_{\bar{\alpha}, i} \varepsilon_{i; \bar{\alpha}} \frac{dz_{i; \bar{\alpha}}}{d\theta_{\ell}}, \quad (86)$$

so we have

$$dz_{i; \delta} = -\eta \sum_{\ell, i', \bar{\alpha}} \varepsilon_{i'; \bar{\alpha}} \frac{dz_{i; \delta}}{d\theta_{\ell}} \frac{dz_{i'; \bar{\alpha}}}{d\theta_{\ell}}. \quad (87)$$

Since we measure the expectation values of operators, our  $\varepsilon$ s are always real. Defining the kernel,

$$K_{\delta, \bar{\alpha}}^{ii'} = \sum_{\ell} \frac{dz_{i; \delta}}{d\theta_{\ell}} \frac{dz_{i'; \bar{\alpha}}}{d\theta_{\ell}}, \quad (88)$$

we have

$$dz_{i; \delta} = -\eta \sum_{\bar{\alpha}, i'} K_{\delta, \bar{\alpha}}^{ii'} \varepsilon_{i'; \bar{\alpha}}. \quad (89)$$

We could also make the joint indices

$$(\delta, i) = \bar{a}, \quad (\bar{\alpha}, i') = \hat{b}, \quad (90)$$

which are running in the space  $\mathcal{D} \times \mathcal{O}(\mathcal{H})$  and  $\mathcal{A} \times \mathcal{O}(\mathcal{H})$  respectively. The notation  $\hat{a}$  is indicating that the data point component belongs to the training set  $\mathcal{A}$ , while the notation  $\bar{a}$  means that the data point component is general in  $\mathcal{D}$ . And we have

$$dz_{\bar{a}} = -\eta \sum_{\hat{b}} K_{\bar{a}\hat{b}} \varepsilon_{\hat{b}}. \quad (91)$$

In general, one could prove the following statement.

**Theorem 1.** *The matrix  $K_{\bar{a}\hat{b}}$  is non-negative and symmetric.*

*Proof.* It is symmetric by definition. Moreover, we consider an arbitrary vector  $f_{\bar{a}}$ . We have

$$\begin{aligned} \sum_{\bar{a}, \hat{b}} f_{\bar{a}} K_{\bar{a}\hat{b}} f_{\hat{b}} &= \sum_{\delta, \delta', i, i'} f_{\delta, i} K_{\delta, \delta'}^{ii'} f_{\delta', i'} = \\ &= \sum_{\delta, \delta', i, i', \ell} f_{\delta, i} f_{\delta', i'} \frac{dz_{i'; \delta'}}{d\theta_{\ell}} \frac{dz_{i; \delta}}{d\theta_{\ell}} = \sum_{\ell} \left( \sum_{\delta, i} f_{\delta, i} \frac{dz_{i; \delta}}{d\theta_{\ell}} \right)^2 \geq 0. \end{aligned} \quad (92)$$

□

Thus, the matrix  $K$  is a proper version of the positive semi-definite symmetric (PDS) kernel [5] in the sense of the classical learning theory.

Now, putting the variational ansatz inside the kernel, we get

$$\begin{aligned} \frac{dz_{i; \delta}}{d\theta_{\ell}} &= \left\langle \phi(\mathbf{x}_{\delta}) \left| U_{+, \ell}^{\dagger} \left[ X_{\ell}, U_{\ell}^{\dagger} W_{\ell}^{\dagger} U_{-, \ell}^{\dagger} O_i U_{-, \ell} W_{\ell} U_{\ell} \right] U_{+, \ell} \right| \phi(\mathbf{x}_{\delta}) \right\rangle, \\ \frac{dz_{i'; \bar{\alpha}}}{d\theta_{\ell}} &= \left\langle \phi(\mathbf{x}_{\bar{\alpha}}) \left| U_{+, \ell}^{\dagger} \left[ X_{\ell}, U_{\ell}^{\dagger} W_{\ell}^{\dagger} U_{-, \ell}^{\dagger} O_{i'} U_{-, \ell} W_{\ell} U_{\ell} \right] U_{+, \ell} \right| \phi(\mathbf{x}_{\bar{\alpha}}) \right\rangle, \end{aligned} \quad (93)$$

So we have

$$K_{\delta, \bar{\alpha}}^{ii'} = \sum_{\ell} \frac{dz_{i; \delta}}{d\theta_{\ell}} \frac{dz_{i'; \bar{\alpha}}}{d\theta_{\ell}} = - \sum_{\ell} \left( \left\langle \phi(\mathbf{x}_{\delta}) \left| U_{+, \ell}^{\dagger} \left[ X_{\ell}, U_{\ell}^{\dagger} W_{\ell}^{\dagger} U_{-, \ell}^{\dagger} O_i U_{-, \ell} W_{\ell} U_{\ell} \right] U_{+, \ell} \right| \phi(\mathbf{x}_{\delta}) \right\rangle \times \left\langle \phi(\mathbf{x}_{\bar{\alpha}}) \left| U_{+, \ell}^{\dagger} \left[ X_{\ell}, U_{\ell}^{\dagger} W_{\ell}^{\dagger} U_{-, \ell}^{\dagger} O_{i'} U_{-, \ell} W_{\ell} U_{\ell} \right] U_{+, \ell} \right| \phi(\mathbf{x}_{\bar{\alpha}}) \right\rangle \right). \quad (94)$$

### B. No representation learning

The statement of *no representation learning* corresponds to the limit where all the change of variational angles are sufficiently close to zero. In this case, the QNTK becomes frozen (static), similar to the optimization problem. With the variational angle redefined, and the frozen QNTK limit, we have,

$$K_{\delta, \bar{\alpha}}^{ii'} = -\delta^2 \sum_{\ell} \left( \left\langle \phi(\mathbf{x}_{\delta}) \left| W_{+, \ell}^{\dagger} \left[ X_{\ell}, W_{\ell}^{\dagger} W_{-, \ell}^{\dagger} O_i W_{-, \ell} W_{\ell} \right] W_{+, \ell} \right| \phi(\mathbf{x}_{\delta}) \right\rangle \times \right. \\ \left. \left\langle \phi(\mathbf{x}_{\bar{\alpha}}) \left| W_{+, \ell}^{\dagger} \left[ X_{\ell}, W_{\ell}^{\dagger} W_{-, \ell}^{\dagger} O_{i'} W_{-, \ell} W_{\ell} \right] W_{+, \ell} \right| \phi(\mathbf{x}_{\bar{\alpha}}) \right\rangle \right). \quad (95)$$

Here,  $\delta$  is the factor we are used to redefine the variational angles  $\theta$  by  $\theta^* + \delta \times \varphi$ . Moreover, one could exactly solve the gradient descent dynamics. We start from the equation

$$d\bar{\varepsilon}_{\hat{a}} = -\eta \sum_{\hat{b}} K_{\hat{a}\hat{b}} \varepsilon_{\hat{b}}. \quad (96)$$

Thus we have

$$\varepsilon_{\hat{a}_1}(t) = \sum_{\hat{a}_2} U_{\hat{a}_1 \hat{a}_2}(t) \varepsilon_{\hat{a}_2}(0), \quad (97)$$

where

$$U_{\hat{a}_1 \hat{a}_2}(t) = \left[ (1 - \eta K)^t \right]_{\hat{a}_1 \hat{a}_2}. \quad (98)$$

One can compute the convergence time as

$$\tau_c = \left\| -\log(1 - \eta K) \right\| \approx \eta \left\| K_{\delta, \bar{\alpha}}^{ii'} \right\|. \quad (99)$$

And we could compute the prediction similarly by

$$\begin{aligned} z_{\bar{a}}(\infty) &= z_{\bar{a}}(0) - \eta \sum_{\hat{b}} K_{\bar{a}\hat{b}} \sum_{t=0}^{+\infty} \varepsilon_{\hat{b}}(t) \\ &= z_{\bar{a}}(0) - \eta \sum_{\hat{b}} K_{\bar{a}\hat{b}} \left( \sum_{t=0, \hat{b}_0}^{+\infty} \left[ (1 - \eta K)^t \right]_{\hat{b}\hat{b}_0} \right) \varepsilon_{\hat{b}_0}(0) \\ &= z_{\bar{a}}(0) - \eta \sum_{\hat{b}} K_{\bar{a}\hat{b}} \sum_{\hat{b}_0} \left[ (1 - (1 - \eta K))^{-1} \right]_{\hat{b}\hat{b}_0} \varepsilon_{\hat{b}_0}(0) \\ &= z_{\bar{a}}(0) - \sum_{\hat{b}} K_{\bar{a}\hat{b}} \sum_{\hat{b}_0} [K^{i, -1}]_{\hat{b}\hat{b}_0} \varepsilon_{\hat{b}_0}(0). \end{aligned} \quad (100)$$

Moreover, we could compute the prediction by firstly defining the kernel inverse. We define

$$\sum_{\hat{a} \in \mathcal{A} \times \mathcal{O}(\mathcal{H})} \tilde{K}^{\hat{a}_1 \hat{a}_2} \tilde{K}_{\hat{a}_2 \hat{a}_3} = \delta_{\hat{a}_3}^{\hat{a}_1}, \quad (101)$$

so we get

$$z_{\bar{a}}(\infty) = z_{\bar{a}}(0) - \sum_{\hat{a}_1, \hat{a}_2} \tilde{K}^{\hat{a}_1 \hat{a}_2} K_{\bar{a}\hat{a}_1} \varepsilon_{\hat{a}_2}(0). \quad (102)$$

### C. Representation learning

Now we start to develop our quantum representation learning theory at the dQNTK order. We make a quadratic expansion:

$$\begin{aligned}
\varepsilon_{i;\bar{\alpha}} &= \langle \phi(\mathbf{x}_{\bar{\alpha}}) | U^\dagger O_i U | \phi(\mathbf{x}_{\bar{\alpha}}) \rangle - y_{i;\bar{\alpha}} \\
&\approx \left\langle \phi(\mathbf{x}_{\bar{\alpha}}) \left| \left( \prod_{\ell'=L}^1 W_{\ell'}^\dagger \right) O_i \left( \prod_{\ell=1}^L W_\ell \right) \right| \phi(\mathbf{x}_{\bar{\alpha}}) \right\rangle - y_{i;\bar{\alpha}} \\
&- i\delta \sum_\ell \varphi_\ell \left\langle \phi(\mathbf{x}_{\bar{\alpha}}) \left| W_{+, \ell}^\dagger \left[ X_\ell, W_\ell^\dagger W_{-, \ell}^\dagger O_i W_{-, \ell} W_\ell \right] W_{+, \ell} \right| \phi(\mathbf{x}_{\bar{\alpha}}) \right\rangle - \frac{\delta^2}{2} \sum_{\ell_1, \ell_2} \varphi_{\ell_1} \varphi_{\ell_2} \times \\
&\left\{ \begin{aligned} &\left\langle \psi(\mathbf{x}_{\bar{\alpha}}) \left| W_{+, \ell_1}^\dagger \left[ X_{\ell_1}, Q_{\ell_1, \ell_2}^\dagger \left[ X_{\ell_2}, W_{\ell_2}^\dagger W_{-, \ell_2}^\dagger O_i W_{-, \ell_2} W_{\ell_2} \right] W_{\ell_2, \ell_1} W_{\ell_1} \right] W_{+, \ell_1} \right| \psi(\mathbf{x}_{\bar{\alpha}}) \right\rangle : \ell_1 \geq \ell_2 \\ &\left\langle \psi(\mathbf{x}_{\bar{\alpha}}) \left| W_{+, \ell_2}^\dagger \left[ X_{\ell_2}, Q_{\ell_2, \ell_1}^\dagger \left[ X_{\ell_1}, W_{\ell_1}^\dagger W_{-, \ell_1}^\dagger O_i W_{-, \ell_1} W_{\ell_1} \right] W_{\ell_1, \ell_2} W_{\ell_2} \right] W_{+, \ell_2} \right| \psi(\mathbf{x}_{\bar{\alpha}}) \right\rangle : \ell_1 < \ell_2 \end{aligned} \right. . \quad (103)
\end{aligned}$$

At the dQNTK order, we have,

$$\dot{z}_{i;\delta} = \sum_\ell \frac{dz_{i;\delta}}{d\varphi_\ell} \dot{\varphi}_\ell - \frac{1}{2} \delta^2 \sum_{\ell_1, \ell_2} \dot{\varphi}_{\ell_1} \dot{\varphi}_{\ell_2} G_{\ell_1, \ell_2}^{\delta, i}, \quad (104)$$

where we define

$$\begin{aligned}
G_{\ell_1, \ell_2}^{\delta, i} &\equiv G_{\ell_1, \ell_2}(\phi(\mathbf{x}_\delta), O_i) \\
&= \left( \begin{aligned} &\left\langle \phi(\mathbf{x}_\delta) \left| W_{+, \ell_1}^\dagger \left[ X_{\ell_1}, Q_{\ell_1, \ell_2}^\dagger \left[ X_{\ell_2}, W_{\ell_2}^\dagger W_{-, \ell_2}^\dagger O_i W_{-, \ell_2} W_{\ell_2} \right] W_{\ell_2, \ell_1} W_{\ell_1} \right] W_{+, \ell_1} \right| \phi(\mathbf{x}_\delta) \right\rangle : \ell_1 \geq \ell_2 \\ &\left\langle \phi(\mathbf{x}_\delta) \left| W_{+, \ell_2}^\dagger \left[ X_{\ell_2}, Q_{\ell_2, \ell_1}^\dagger \left[ X_{\ell_1}, W_{\ell_1}^\dagger W_{-, \ell_1}^\dagger O_i W_{-, \ell_1} W_{\ell_1} \right] W_{\ell_1, \ell_2} W_{\ell_2} \right] W_{+, \ell_2} \right| \phi(\mathbf{x}_\delta) \right\rangle : \ell_1 < \ell_2 \end{aligned} \right), \\
\Theta_\ell^{\delta, i} &\equiv \Theta_\ell(\phi(\mathbf{x}_\delta), O_i) = \left\langle \phi(\mathbf{x}_\delta) \left| W_{+, \ell}^\dagger \left[ X_\ell, W_\ell^\dagger W_{-, \ell}^\dagger O_i W_{-, \ell} W_\ell \right] W_{+, \ell} \right| \phi(\mathbf{x}_\delta) \right\rangle . \quad (105)
\end{aligned}$$

Moreover, we want to define

$$\Theta_\ell^{\delta, i} \equiv \Theta_\ell(\phi(\mathbf{x}_\delta), O_i) = \left\langle \phi(\mathbf{x}_\delta) \left| W_{+, \ell}^\dagger \left[ X_\ell, W_\ell^\dagger W_{-, \ell}^\dagger O_i W_{-, \ell} W_\ell \right] W_{+, \ell} \right| \phi(\mathbf{x}_\delta) \right\rangle . \quad (106)$$

We will now compute the effective kernel. We have,

$$\begin{aligned}
\frac{dz_{i;\delta}}{d\varphi_\ell} &= -i\delta \Theta_\ell^{\delta, i} - \delta^2 \sum_{\ell'} \varphi_{\ell'} G_{\ell', \ell}^{\delta, i}, \\
\frac{dz_{i';\bar{\alpha}}}{d\varphi_\ell} &= -i\delta \Theta_\ell^{\bar{\alpha}, i'} - \delta^2 \sum_{\ell'} \varphi_{\ell'} G_{\ell', \ell}^{\bar{\alpha}, i'}. \quad (107)
\end{aligned}$$

Thus, we define

$$K_{\delta, \bar{\alpha}}^{E, ii'} = \sum_\ell \frac{dz_{i;\delta}}{d\varphi_\ell} \frac{dz_{i';\bar{\alpha}}}{d\varphi_\ell}, \quad (108)$$

and

$$K_{\delta, \bar{\alpha}}^{E, ii'}(0) = K_{\delta, \bar{\alpha}}^{ii'} + K_{\delta, \bar{\alpha}}^{\Delta, ii'}. \quad (109)$$

We find

$$K_{\delta, \bar{\alpha}}^{\Delta, ii'} = i\delta^3 \sum_{\ell, \ell'} \varphi_{\ell'}(0) G_{\ell', \ell}^{\delta, i} \Theta_\ell^{\bar{\alpha}, i'} + i\delta^3 \sum_{\ell, \ell'} \varphi_{\ell'}(0) G_{\ell', \ell}^{\bar{\alpha}, i'} \Theta_\ell^{\delta, i}. \quad (110)$$

Now we could write down the prediction on the training set. We have

$$\varepsilon_{\hat{a}}(t) = \varepsilon_{\hat{a}}^F(t) + \varepsilon_{\hat{a}}^I(t), \quad (111)$$

where

$$\begin{aligned}\varepsilon_{\hat{a}}^F(t) &= \sum_{\hat{a}_1} U_{\hat{a}\hat{a}_1}(t) \varepsilon_{\hat{a}_1}(0), \\ U_{\hat{a}_1\hat{a}_2}(t) &= \left[ (1 - \eta K)^t \right]_{\hat{a}_1\hat{a}_2},\end{aligned}\quad (112)$$

and

$$\varepsilon_{\hat{a}}^I(t) = \left( -\eta \sum_{s=0}^{t-1} (1 - \eta K)^{t-1-s} K^\Delta (1 - \eta K)^s \varepsilon(0) \right)_{\hat{a}}. \quad (113)$$

It is the compact matrix product form in the space  $\mathcal{A} \times \mathcal{O}(\mathcal{H})$ . Moreover, we have

$$\|\varepsilon^I(t)\| \leq \eta t \|1 - \eta K\|^{t-1} \|K^\Delta\| \|\varepsilon(0)\|. \quad (114)$$

Finally, we discuss the asymptotic convergence regime. We could compute the quantum meta-kernel. We notice that

$$dz_{i;\delta} = -\eta \sum_{\ell, i', \tilde{\alpha}} \frac{dz_{i;\delta}}{d\varphi_\ell} \frac{dz_{i';\tilde{\alpha}}}{d\varphi_\ell} \varepsilon_{i;\tilde{\alpha}} + \eta^2 \sum_{\ell_1, \ell_2, \tilde{\alpha}_1, i_1, \tilde{\alpha}_2, i_2} \frac{d^2 z_{i;\delta}}{d\varphi_{\ell_1} d\varphi_{\ell_2}} \frac{dz_{i_1;\tilde{\alpha}_1}}{d\varphi_{\ell_1}} \frac{dz_{i_2;\tilde{\alpha}_2}}{d\varphi_{\ell_2}} \varepsilon_{i_1;\tilde{\alpha}_1} \varepsilon_{i_2;\tilde{\alpha}_2}. \quad (115)$$

So we define

$$\mu_{\delta_0 \tilde{\alpha}_1 \tilde{\alpha}_2}^{i_1 i_2} = \mu_{\bar{a} \hat{a}_1 \hat{a}_2} = \sum_{\ell_1, \ell_2} \frac{d^2 z_{i;\delta}}{d\varphi_{\ell_1} d\varphi_{\ell_2}} \left( \frac{dz_{i_1;\tilde{\alpha}_1}}{d\varphi_{\ell_1}} \frac{dz_{i_2;\tilde{\alpha}_2}}{d\varphi_{\ell_2}} \right) \Big|_{\varphi=0}, \quad (116)$$

in the leading order. Note that it is natural to extend the definition to more general inputs

$$\mu_{\delta_0 \delta_1 \delta_2}^{i_0 i_1 i_2} = \mu_{\bar{a}_0 \bar{a}_1 \bar{a}_2} = \sum_{\ell_1, \ell_2} \frac{d^2 z_{i_0;\delta_0}}{d\varphi_{\ell_1} d\varphi_{\ell_2}} \left( \frac{dz_{i_1;\delta_1}}{d\varphi_{\ell_1}} \frac{dz_{i_2;\delta_2}}{d\varphi_{\ell_2}} \right) \Big|_{\varphi=0}. \quad (117)$$

So the asymptotic convergence is given by

$$\begin{aligned}z_{\bar{a}}(\infty) &= z_{\bar{a}}(0) - \sum_{\hat{a}_1, \hat{a}_2} K_{\bar{a}\hat{a}_1} \tilde{K}^{\hat{a}_1 \hat{a}_2} \varepsilon_{\hat{a}_2}(0) \\ &+ \sum_{\hat{a}_1, \hat{a}_2, \hat{a}_3, \hat{a}_4} \left[ \mu_{\hat{a}_1 \bar{a} \hat{a}_2} - \sum_{\hat{a}_5, \hat{a}_6} K_{\bar{a}\hat{a}_5} \tilde{K}^{\hat{a}_5 \hat{a}_6} \mu_{\hat{a}_1 \hat{a}_6 \hat{a}_2} \right] Z_A^{\hat{a}_1 \hat{a}_2 \hat{a}_3 \hat{a}_4} \varepsilon_{\hat{a}_3}(0) \varepsilon_{\hat{a}_4}(0) \\ &+ \sum_{\hat{a}_1, \hat{a}_2, \hat{a}_3, \hat{a}_4} \left[ \mu_{\bar{a} \hat{a}_1 \hat{a}_2} - \sum_{\hat{a}_5, \hat{a}_6} K_{\bar{a}\hat{a}_5} \tilde{K}^{\hat{a}_5 \hat{a}_6} \mu_{\hat{a}_6 \hat{a}_1 \hat{a}_2} \right] Z_B^{\hat{a}_1 \hat{a}_2 \hat{a}_3 \hat{a}_4} \varepsilon_{\hat{a}_3}(0) \varepsilon_{\hat{a}_4}(0),\end{aligned}\quad (118)$$

where the algorithm projector  $Z_{A,B}$ s are

$$\begin{aligned}Z_A^{\hat{a}_1 \hat{a}_2 \hat{a}_3 \hat{a}_4} &\equiv \tilde{K}^{\hat{a}_1 \hat{a}_3} \tilde{K}^{\hat{a}_2 \hat{a}_4} - \sum_{\hat{a}_5} \tilde{K}^{\hat{a}_2 \hat{a}_5} X_{\parallel}^{\hat{a}_1 \hat{a}_5 \hat{a}_3 \hat{a}_4}, \\ Z_B^{\hat{a}_1 \hat{a}_2 \hat{a}_3 \hat{a}_4} &\equiv \tilde{K}^{\hat{a}_1 \hat{a}_3} \tilde{K}^{\hat{a}_2 \hat{a}_4} - \sum_{\hat{a}_5} \tilde{K}^{\hat{a}_2 \hat{a}_5} X_{\parallel}^{\hat{a}_1 \hat{a}_5 \hat{a}_3 \hat{a}_4} + \frac{\eta}{2} X_{\parallel}^{\hat{a}_1 \hat{a}_2 \hat{a}_3 \hat{a}_4},\end{aligned}\quad (119)$$

and

$$X_{\parallel}^{\hat{a}_1 \hat{a}_2 \hat{a}_3 \hat{a}_4} = \sum_{s=0}^{\infty} [(1 - \eta K)^s]_{\hat{a}_1 \hat{a}_3} [(1 - \eta K)^s]_{\hat{a}_2 \hat{a}_4}, \quad (120)$$

which is defined implicitly as

$$\delta_{\hat{a}_5}^{\hat{a}_1} \delta_{\hat{a}_6}^{\hat{a}_2} = \sum_{\hat{a}_3, \hat{a}_4} X_{\parallel}^{\hat{a}_1 \hat{a}_2 \hat{a}_3 \hat{a}_4} \left( \tilde{K}_{\hat{a}_3 \hat{a}_5} \delta_{\hat{a}_4 \hat{a}_6} + \delta_{\hat{a}_3 \hat{a}_5} \tilde{K}_{\hat{a}_4 \hat{a}_6} - \eta \tilde{K}_{\hat{a}_3 \hat{a}_5} \tilde{K}_{\hat{a}_4 \hat{a}_6} \right). \quad (121)$$

Now we write down all components of  $\mu$ . We have

$$\begin{aligned}\frac{dz_{i_1;\delta_1}}{d\varphi_{\ell_1}}(\varphi=0) &= -i\delta\Theta_{\ell_1}^{\delta_1,i_1}, \\ \frac{dz_{i_2;\delta_2}}{d\varphi_{\ell_2}}(\varphi=0) &= -i\delta\Theta_{\ell_2}^{\delta_2,i_2}, \\ \frac{d^2z_{i_0;\delta_0}}{d\varphi_{\ell_1}d\varphi_{\ell_2}} &= -\delta^2G_{\ell_1,\ell_2}^{\delta_0,i_0}.\end{aligned}\quad (122)$$

So

$$\mu_{\delta_0\delta_1\delta_2}^{i_0i_1i_2} = \sum_{\ell_1,\ell_2} \frac{d^2z_{i_0;\delta_0}}{d\varphi_{\ell_1}d\varphi_{\ell_2}} \left( \frac{dz_{i_1;\delta_1}}{d\varphi_{\ell_1}} \frac{dz_{i_2;\delta_2}}{d\varphi_{\ell_2}} \right) \Big|_{\varphi=0} = \delta^4 \sum_{\ell_1,\ell_2} \Theta_{\ell_1}^{\delta_1,i_1} \Theta_{\ell_2}^{\delta_2,i_2} G_{\ell_1,\ell_2}^{\delta_0,i_0}.\quad (123)$$

#### IV. DETAILS ON THE QUANTUM MACHINE LEARNING: AMPLITUDE ENCODING

##### A. General setup

Here, we consider alternative quantum machine learning models with the amplitude encoding. In this case, we naturally extend the NTK formalism to the complex version. We consider the variational setup where

$$z_{i;\delta} \equiv z_i(\theta, \mathbf{x}_\delta) = \left\langle i \left| \prod_{\ell=1}^L W_\ell U_\ell \right| \phi(\mathbf{x}_\delta) \right\rangle.\quad (124)$$

One could take the derivative,

$$\bar{d}z_{i;\delta} = \sum_{\ell} \frac{dz_{i;\delta}}{d\theta_{\ell}} d\theta_{\ell},\quad (125)$$

The loss is

$$L_{\mathcal{A}}(\theta) = \frac{1}{2} \sum_{\bar{\alpha},i} |y_{i;\bar{\alpha}} - z_{i;\bar{\alpha}}|^2 = \frac{1}{2} \sum_{\bar{\alpha},i} (y_{i;\bar{\alpha}} - z_{i;\bar{\alpha}})(y_{i;\bar{\alpha}} - z_{i;\bar{\alpha}})^*.\quad (126)$$

So

$$\frac{dL_{\mathcal{A}}(\theta)}{d\theta_{\ell}} = \frac{1}{2} \sum_{\bar{\alpha},i} \varepsilon_{i;\bar{\alpha}} \frac{dz_{i;\bar{\alpha}}^*}{d\theta_{\ell}} + \frac{1}{2} \sum_{\bar{\alpha},i} \varepsilon_{i;\bar{\alpha}}^* \frac{dz_{i;\bar{\alpha}}}{d\theta_{\ell}} = \text{Re} \sum_{\bar{\alpha},i} \varepsilon_{i;\bar{\alpha}} \frac{dz_{i;\bar{\alpha}}^*}{d\theta_{\ell}}.\quad (127)$$

Note that if we count the number of times that  $U$  appears in the loss function, the amplitude encoding model is a *squareroot* of the operator expectation value model. The gradient descent rule is

$$d\theta_{\ell} = -\eta \frac{d\mathcal{L}_{\mathcal{A}}}{d\theta_{\ell}} = -\frac{\eta}{2} \sum_{\bar{\alpha},i} \varepsilon_{i;\bar{\alpha}} \frac{dz_{i;\bar{\alpha}}^*}{d\theta_{\ell}} - \frac{\eta}{2} \sum_{\bar{\alpha},i} \varepsilon_{i;\bar{\alpha}}^* \frac{dz_{i;\bar{\alpha}}}{d\theta_{\ell}},\quad (128)$$

so we have

$$\begin{aligned}\bar{d}z_{i;\delta} &= -\eta \sum_{\ell,i',\bar{\alpha}} \frac{dz_{i;\delta}}{d\theta_{\ell}} \text{Re} \left( \varepsilon_{i';\bar{\alpha}} \frac{dz_{i';\bar{\alpha}}^*}{d\theta_{\ell}} \right) \\ &= -\frac{\eta}{2} \sum_{\ell,i',\bar{\alpha}} \varepsilon_{i';\bar{\alpha}} \frac{dz_{i;\delta}}{d\theta_{\ell}} \frac{dz_{i';\bar{\alpha}}^*}{d\theta_{\ell}} - \frac{\eta}{2} \sum_{\ell,i',\bar{\alpha}} \varepsilon_{i';\bar{\alpha}}^* \frac{dz_{i;\delta}}{d\theta_{\ell}} \frac{dz_{i';\bar{\alpha}}}{d\theta_{\ell}}.\end{aligned}\quad (129)$$

We notice that the variable  $z$  is complex in general. So we could write,

$$\begin{aligned}\bar{d}z_{i;\delta} &= -\frac{\eta}{2} \sum_{\ell,i',\bar{\alpha}} \varepsilon_{i';\bar{\alpha}} \frac{dz_{i;\delta}}{d\theta_{\ell}} \frac{dz_{i';\bar{\alpha}}^*}{d\theta_{\ell}} - \frac{\eta}{2} \sum_{\ell,i',\bar{\alpha}} \varepsilon_{i';\bar{\alpha}}^* \frac{dz_{i;\delta}}{d\theta_{\ell}} \frac{dz_{i';\bar{\alpha}}}{d\theta_{\ell}}, \\ \bar{d}z_{i;\delta}^* &= -\frac{\eta}{2} \sum_{\ell,i',\bar{\alpha}} \varepsilon_{i';\bar{\alpha}}^* \frac{dz_{i;\delta}^*}{d\theta_{\ell}} \frac{dz_{i';\bar{\alpha}}}{d\theta_{\ell}} - \frac{\eta}{2} \sum_{\ell,i',\bar{\alpha}} \varepsilon_{i';\bar{\alpha}} \frac{dz_{i;\delta}^*}{d\theta_{\ell}} \frac{dz_{i';\bar{\alpha}}^*}{d\theta_{\ell}}.\end{aligned}\quad (130)$$

Defining the kernel

$$\begin{pmatrix} K_{\delta,\tilde{\alpha}}^{+,ii'} & K_{\delta,\tilde{\alpha}}^{-,ii'} \\ K_{\delta,\tilde{\alpha}}^{*,-,ii'} & K_{\delta,\tilde{\alpha}}^{*,+,ii'} \end{pmatrix} = \begin{pmatrix} \sum_{\ell} \frac{dz_{i;\delta}}{d\theta_{\ell}} \frac{dz_{i';\tilde{\alpha}}^*}{d\theta_{\ell}} & \sum_{\ell} \frac{dz_{i;\delta}}{d\theta_{\ell}} \frac{dz_{i';\tilde{\alpha}}}{d\theta_{\ell}} \\ \sum_{\ell} \frac{dz_{i;\delta}^*}{d\theta_{\ell}} \frac{dz_{i';\tilde{\alpha}}}{d\theta_{\ell}} & \sum_{\ell} \frac{dz_{i;\delta}^*}{d\theta_{\ell}} \frac{dz_{i';\tilde{\alpha}}^*}{d\theta_{\ell}} \end{pmatrix}, \quad (131)$$

we have

$$\begin{aligned} \vec{dz}_{i;\delta} &= -\frac{\eta}{2} \sum_{\tilde{\alpha},i'} K_{\delta,\tilde{\alpha}}^{+,ii'} \varepsilon_{i';\tilde{\alpha}} - \frac{\eta}{2} \sum_{\tilde{\alpha},i'} K_{\delta,\tilde{\alpha}}^{-,ii'} \varepsilon_{i';\tilde{\alpha}}^*, \\ \vec{dz}_{i;\delta}^* &= -\frac{\eta}{2} \sum_{\tilde{\alpha},i'} K_{\delta,\tilde{\alpha}}^{*,-,ii'} \varepsilon_{i';\tilde{\alpha}} - \frac{\eta}{2} \sum_{\tilde{\alpha},i'} K_{\delta,\tilde{\alpha}}^{*,+,ii'} \varepsilon_{i';\tilde{\alpha}}^*, \end{aligned} \quad (132)$$

Here we also use  $K^*$  to denote the complex conjugate. We introduce the worldsheet index  $\tilde{\alpha}, \tilde{\beta}$  to denote the two-component of the complex variable  $(\vec{dz}, \vec{dz}^*)$ . So we have [6]

$$\vec{dz}_{i;\delta}^{\tilde{\alpha}} = -\frac{\eta}{2} \sum_{\tilde{\alpha},\tilde{\beta},i'} K_{\delta,\tilde{\alpha}}^{-\tilde{\alpha}\tilde{\beta},ii'} \varepsilon_{i';\tilde{\alpha}}^{\tilde{\beta}}, \quad (133)$$

where

$$\varepsilon_{i;\delta}^{\tilde{\alpha}} = \begin{pmatrix} \varepsilon_{i;\delta} \\ \varepsilon_{i;\delta}^* \end{pmatrix}, \quad K_{\delta,\tilde{\alpha}}^{-\tilde{\alpha}\tilde{\beta},ii'} = \begin{pmatrix} K_{\delta,\tilde{\alpha}}^{+,ii'} & K_{\delta,\tilde{\alpha}}^{-,ii'} \\ K_{\delta,\tilde{\alpha}}^{*,-,ii'} & K_{\delta,\tilde{\alpha}}^{*,+,ii'} \end{pmatrix}. \quad (134)$$

We could also make the joint index

$$(\delta, \tilde{\alpha}, i) = \bar{\mu}, \quad (\tilde{\alpha}, \tilde{\beta}, i') = \hat{\nu}, \quad (135)$$

which is running in the space  $\mathcal{D} \times \mathbb{Z}_2 \times \mathcal{H}$  and  $\mathcal{A} \times \mathbb{Z}_2 \times \mathcal{H}$  respectively. The notation  $\hat{\mu}$  is indicating that the data point component belongs to the training set  $\mathcal{A}$ , while the notation  $\bar{\mu}$  means that the data point component is general in  $\mathcal{D}$ . And we have

$$\vec{dz}_{\bar{\mu}} = -\frac{\eta}{2} \sum_{\hat{\nu}} K_{\bar{\mu}\hat{\nu}} \varepsilon_{\hat{\nu}}. \quad (136)$$

In general, one could prove the following statement.

**Theorem 2.** *The matrix  $K_{\bar{\mu}\hat{\nu}}$  is non-negative and Hermitian.*

*Proof.* One could easily check that the matrix form

$$\begin{aligned} K_{\bar{\mu}\hat{\nu}}^{\dagger} &= \begin{pmatrix} K_{\delta,\tilde{\alpha}}^{+,ii'} & K_{\delta,\tilde{\alpha}}^{-,ii'} \\ K_{\delta,\tilde{\alpha}}^{*,-,ii'} & K_{\delta,\tilde{\alpha}}^{*,+,ii'} \end{pmatrix}^{*,T} = \begin{pmatrix} K_{\delta,\tilde{\alpha}}^{*,-,ii'} & K_{\delta,\tilde{\alpha}}^{*,+,ii'} \\ K_{\delta,\tilde{\alpha}}^{-,ii'} & K_{\delta,\tilde{\alpha}}^{+,ii'} \end{pmatrix}^T \\ &= \begin{pmatrix} (K_{\delta,\tilde{\alpha}}^{*,-,ii'})^T & (K_{\delta,\tilde{\alpha}}^{*,+,ii'})^T \\ (K_{\delta,\tilde{\alpha}}^{-,ii'})^T & (K_{\delta,\tilde{\alpha}}^{+,ii'})^T \end{pmatrix} = \begin{pmatrix} K_{\delta,\tilde{\alpha}}^{+,ii'} & K_{\delta,\tilde{\alpha}}^{-,ii'} \\ K_{\delta,\tilde{\alpha}}^{*,-,ii'} & K_{\delta,\tilde{\alpha}}^{*,+,ii'} \end{pmatrix} = K_{\bar{\mu}\hat{\nu}}, \end{aligned} \quad (137)$$

is Hermitian. Moreover, consider an arbitrary vector  $f_{\bar{\mu}}$ , we have

$$\begin{aligned} \sum_{\bar{\mu},\hat{\nu}} f_{\bar{\mu}}^* K_{\bar{\mu}\hat{\nu}} f_{\hat{\nu}} &= \sum_{\delta,\delta',i,i'} (f_{\delta,i}^* \ f_{\delta,i}) \begin{pmatrix} K_{\delta,\delta'}^{+,ii'} & K_{\delta,\delta'}^{-,ii'} \\ K_{\delta,\delta'}^{*,-,ii'} & K_{\delta,\delta'}^{*,+,ii'} \end{pmatrix} \begin{pmatrix} f_{\delta',i'} \\ f_{\delta',i'}^* \end{pmatrix} \\ &= \sum_{\delta,\delta',i,i'} \left( f_{\delta,i}^* f_{\delta',i'} K_{\delta,\delta'}^{+,ii'} + f_{\delta,i} f_{\delta',i'} K_{\delta,\delta'}^{-,ii'} + f_{\delta,i}^* f_{\delta',i'}^* K_{\delta,\delta'}^{*,-,ii'} + f_{\delta,i} f_{\delta',i'} K_{\delta,\delta'}^{*,+,ii'} \right) \\ &= \sum_{\delta,\delta',i,i',\ell} \left( f_{\delta,i}^* f_{\delta',i'} \frac{dz_{i;\delta}}{d\theta_{\ell}} \frac{dz_{i';\delta'}^*}{d\theta_{\ell}} + f_{\delta,i} f_{\delta',i'} \frac{dz_{i;\delta}^*}{d\theta_{\ell}} \frac{dz_{i';\delta'}}{d\theta_{\ell}} \right. \\ &\quad \left. + f_{\delta,i}^* f_{\delta',i'}^* \frac{dz_{i;\delta}}{d\theta_{\ell}} \frac{dz_{i';\delta'}}{d\theta_{\ell}} + f_{\delta,i} f_{\delta',i'} \frac{dz_{i;\delta}^*}{d\theta_{\ell}} \frac{dz_{i';\delta'}^*}{d\theta_{\ell}} \right) \\ &= \sum_{\delta,\delta',i,i',\ell} \left( f_{\delta,i} \frac{dz_{i;\delta}^*}{d\theta_{\ell}} + f_{\delta,i}^* \frac{dz_{i;\delta}}{d\theta_{\ell}} \right) \left( f_{\delta',i'} \frac{dz_{i';\delta'}}{d\theta_{\ell}} + f_{\delta',i'}^* \frac{dz_{i';\delta'}^*}{d\theta_{\ell}} \right) \\ &= \sum_{\ell} \left( \sum_{\delta,i} f_{\delta,i} \frac{dz_{i;\delta}^*}{d\theta_{\ell}} + f_{\delta,i}^* \frac{dz_{i;\delta}}{d\theta_{\ell}} \right)^2 = 4 \sum_{\ell} \left( \text{Re} \sum_{\delta,i} f_{\delta,i} \frac{dz_{i;\delta}^*}{d\theta_{\ell}} \right)^2 \geq 0. \end{aligned} \quad (138)$$

□

Thus, the matrix  $K$  is a complexified version of the positive semi-definite symmetric (PDS) kernel in the sense of the kernel method of statistical learning theory (see, for instance, [5]), but running during the training dynamics. Moreover, our work provides a complexified version of the NTK theory that could be useful for machine learning itself.

Now, putting the variational ansatz inside the kernel, we get

$$\frac{dz_{i;\alpha}}{d\theta_\ell} = \frac{d}{d\theta_\ell} \left\langle i \left| \prod_{\ell'=1}^L W_{\ell'} U_{\ell'} \right| \phi(\mathbf{x}_\alpha) \right\rangle = i \left\langle i \left| \prod_{\ell'=1}^{\ell-1} W_{\ell'} U_{\ell'} X_\ell \prod_{\ell''=\ell}^L W_{\ell''} U_{\ell''} \right| \phi(\mathbf{x}_\alpha) \right\rangle, \quad (139)$$

$$\frac{dz_{i;\alpha}^*}{d\theta_\ell} = \frac{d}{d\theta_\ell} \left\langle \phi(\mathbf{x}_\alpha) \left| \prod_{\ell'=L}^1 U_{\ell'}^\dagger W_{\ell'}^\dagger \right| i \right\rangle = -i \left\langle \phi(\mathbf{x}_\alpha) \left| \prod_{\ell'=L}^{\ell+1} U_{\ell'}^\dagger W_{\ell'}^\dagger X_\ell \prod_{\ell''=1}^{\ell} U_{\ell''}^\dagger W_{\ell''}^\dagger \right| i \right\rangle. \quad (140)$$

And we have

$$\begin{aligned} \frac{dz_{i;\alpha}}{d\theta_\ell} &= i \langle i | U_{-, \ell} X_\ell W_\ell U_\ell U_{+, \ell} | \phi(\mathbf{x}_\alpha) \rangle, \\ \frac{dz_{i;\alpha}^*}{d\theta_\ell} &= -i \langle \phi(\mathbf{x}_\alpha) | U_{+, \ell}^\dagger U_\ell^\dagger W_\ell^\dagger X_\ell U_{-, \ell}^\dagger | i \rangle. \end{aligned} \quad (141)$$

So we have

$$\begin{aligned} \sum_\ell \frac{dz_{i;\delta}}{d\theta_\ell} \frac{dz_{i';\bar{\alpha}}^*}{d\theta_\ell} &= \sum_\ell \langle i | U_{-, \ell} X_\ell W_\ell U_\ell U_{+, \ell} | \phi(\mathbf{x}_\delta) \rangle \langle \phi(\mathbf{x}_{\bar{\alpha}}) | U_{+, \ell}^\dagger U_\ell^\dagger W_\ell^\dagger X_\ell U_{-, \ell}^\dagger | i' \rangle, \\ \sum_\ell \frac{dz_{i;\delta}}{d\theta_\ell} \frac{dz_{i';\bar{\alpha}}}{d\theta_\ell} &= - \sum_\ell \langle i | U_{-, \ell} X_\ell W_\ell U_\ell U_{+, \ell} | \phi(\mathbf{x}_\delta) \rangle \langle i' | U_{-, \ell} X_\ell W_\ell U_\ell U_{+, \ell} | \phi(\mathbf{x}_{\bar{\alpha}}) \rangle, \end{aligned} \quad (142)$$

and their conjugates. Now, we define some notations. We define the *feature S-matrix*

$$\rho_{\delta\bar{\alpha}} = |\phi(\mathbf{x}_\delta)\rangle \langle \phi(\mathbf{x}_{\bar{\alpha}})|, \quad (143)$$

and the *feature projector*

$$\Lambda_{i';\delta} = |\phi(\mathbf{x}_\delta)\rangle \langle i'|. \quad (144)$$

We have

$$\begin{aligned} K_{\delta,\bar{\alpha}}^{+,ii'} &= \sum_\ell \langle i | U_{-, \ell} X_\ell W_\ell U_\ell U_{+, \ell} \rho_{\delta\bar{\alpha}} U_{+, \ell}^\dagger U_\ell^\dagger W_\ell^\dagger X_\ell U_{-, \ell}^\dagger | i' \rangle, \\ K_{\delta,\bar{\alpha}}^{-,ii'} &= - \sum_\ell \langle i' | U_{-, \ell} X_\ell W_\ell U_\ell U_{+, \ell} \Lambda_{i';\delta} U_{-, \ell} X_\ell W_\ell U_\ell U_{+, \ell} | \phi(\mathbf{x}_\delta) \rangle. \end{aligned} \quad (145)$$

## B. No representation learning

In the frozen QNTK limit with the variational angle redefined, one could compute the expressions of the kernel as

$$\begin{aligned} K_{\delta,\bar{\alpha}}^{+,ii'} &= \delta^2 \sum_\ell \langle i | W_{-, \ell} X_\ell W_\ell W_{+, \ell} \rho_{\delta\bar{\alpha}} W_{+, \ell}^\dagger W_\ell^\dagger X_\ell W_{-, \ell}^\dagger | i' \rangle, \\ K_{\delta,\bar{\alpha}}^{-,ii'} &= -\delta^2 \sum_\ell \langle i' | W_{-, \ell} X_\ell W_\ell U_{+, \ell} \Lambda_{i';\delta} W_{-, \ell} X_\ell W_\ell W_{+, \ell} | \phi(\mathbf{x}_\delta) \rangle. \end{aligned} \quad (146)$$

Here,  $\delta$  is the factor we are used to redefine the variational angles  $\theta$  by  $\theta^* + \delta \times \varphi$ , and we define

$$W_{-, \ell} \equiv \prod_{\ell'=1}^{\ell-1} W_{\ell'}, \quad W_{+, \ell} \equiv \prod_{\ell'=\ell+1}^L W_{\ell'}. \quad (147)$$

Moreover, one could exactly solve the gradient descent dynamics. We start from the equation

$$d\bar{\varepsilon}_{\hat{\mu}} = -\frac{\eta}{2} \sum_{\hat{\nu}} K_{\hat{\mu}\hat{\nu}} \varepsilon_{\hat{\nu}}. \quad (148)$$

Thus we have

$$\varepsilon_{\hat{\mu}_1}(t) = \sum_{\hat{\mu}_2} U_{\hat{\mu}_1 \hat{\mu}_2}(t) \varepsilon_{\hat{\mu}_2}(0), \quad (149)$$

where

$$U_{\hat{\mu}_1 \hat{\mu}_2}(t) = \left[ \left(1 - \frac{\eta}{2} K\right)^t \right]_{\hat{\mu}_1 \hat{\mu}_2}. \quad (150)$$

One can compute the convergence time as

$$\tau_c = \left\| -\log \left(1 - \frac{\eta}{2} K\right) \right\| \approx \frac{\eta}{2} \left\| \begin{pmatrix} K_{\delta, \bar{\alpha}}^{+, ii'} & K_{\delta, \bar{\alpha}}^{-, ii'} \\ K_{\delta, \bar{\alpha}}^{*, -, ii'} & K_{\delta, \bar{\alpha}}^{*, +, ii'} \end{pmatrix} \right\|. \quad (151)$$

And we could compute the prediction similarly by

$$\begin{aligned} z_{\bar{\mu}}(\infty) &= z_{\bar{\mu}}(0) - \frac{\eta}{2} \sum_{\hat{\nu}} K_{\bar{\mu} \hat{\nu}} \sum_{t=0}^{+\infty} \varepsilon_{\hat{\nu}}(t) \\ &= z_{\bar{\mu}}(0) - \frac{\eta}{2} \sum_{\hat{\nu}} K_{\bar{\mu} \hat{\nu}} \left( \sum_{t=0, \hat{\nu}_0}^{+\infty} \left[ \left(1 - \frac{\eta}{2} K\right)^t \right]_{\hat{\nu} \hat{\nu}_0} \right) \varepsilon_{\hat{\nu}_0}(0) \\ &= z_{\bar{\mu}}(0) - \frac{\eta}{2} \sum_{\hat{\nu}} K_{\bar{\mu} \hat{\nu}} \sum_{\hat{\nu}_0} \left[ \left(1 - \left(1 - \frac{\eta}{2} K\right)\right)^{-1} \right]_{\hat{\nu} \hat{\nu}_0} \varepsilon_{\hat{\nu}_0}(0) \\ &= z_{\bar{\mu}}(0) - \sum_{\hat{\nu}} K_{\bar{\mu} \hat{\nu}} \sum_{\hat{\nu}_0} [K^{i, -1}]_{\hat{\nu} \hat{\nu}_0} \varepsilon_{\hat{\nu}_0}(0). \end{aligned} \quad (152)$$

Moreover, we could compute the prediction by firstly defining the kernel inverse. We define

$$\sum_{\hat{\mu} \in \mathcal{A} \times \mathbb{Z}_2 \times \mathcal{H}} \tilde{K}^{\hat{\mu}_1 \hat{\mu}_2} \tilde{K}^{\hat{\mu}_2 \hat{\mu}_3} = \delta_{\hat{\mu}_3}^{\hat{\mu}_1}, \quad (153)$$

so we get

$$z_{\bar{\mu}}(\infty) = z_{\bar{\mu}}(0) - \sum_{\hat{\mu}_1, \hat{\mu}_2} \tilde{K}^{\hat{\mu}_1 \hat{\mu}_2} K_{\bar{\mu} \hat{\mu}_1} \varepsilon_{\hat{\mu}_2}(0). \quad (154)$$

### C. Representation learning

Now we start to develop our quantum representation learning theory at the dQNTK order. We make a quadratic expansion:

$$\begin{aligned} U(\theta) &\rightarrow \prod_{\ell} W_{\ell} \exp(i\delta \varphi_{\ell} X_{\ell}) = \prod_{\ell} W_{\ell} + i\delta \sum_{\ell} \varphi_{\ell} W_{\ell, -} W_{\ell} X_{\ell} W_{\ell, +} \\ &- \frac{1}{2} \delta^2 \sum_{\ell_1, \ell_2} \varphi_{\ell_1} \varphi_{\ell_2} \begin{cases} W_{\ell_1, -} W_{\ell_1} X_{\ell_1} W_{\ell_1, \ell_2} W_{\ell_2} X_{\ell_2} W_{\ell_2, +} : \ell_1 \leq \ell_2 \\ W_{\ell_2, -} W_{\ell_2} X_{\ell_2} W_{\ell_2, \ell_1} W_{\ell_1} X_{\ell_1} W_{\ell_1, +} : \ell_1 > \ell_2 \end{cases}. \end{aligned} \quad (155)$$

Now, we could call

$$\begin{aligned} G_{\ell}^{(1)} &\equiv W_{\ell} X_{\ell}, \\ G_{\ell_1 \ell_2}^{(2)} &\equiv \begin{cases} W_{\ell_1, -} W_{\ell_1} X_{\ell_1} W_{\ell_1, \ell_2} W_{\ell_2} X_{\ell_2} W_{\ell_2, +} : \ell_1 \leq \ell_2 \\ W_{\ell_2, -} W_{\ell_2} X_{\ell_2} W_{\ell_2, \ell_1} W_{\ell_1} X_{\ell_1} W_{\ell_1, +} : \ell_1 > \ell_2 \end{cases}. \end{aligned} \quad (156)$$

So the model will look like

$$z_{j; \delta} = \sum_i \left( \prod_{\ell} W_{\ell} \right)_{ji} \phi_i(\mathbf{x}_{\delta}) + i\delta \sum_{i, \ell} \varphi_{\ell} G_{\ell, ji}^{(1)} \phi_i(\mathbf{x}_{\delta}) - \frac{1}{2} \delta^2 \sum_{i, \ell_1, \ell_2} \varphi_{\ell_1} \varphi_{\ell_2} G_{\ell_1 \ell_2, ji}^{(2)} \phi_i(\mathbf{x}_{\delta}). \quad (157)$$

So we could write down the derivatives

$$\frac{dz_{j;\delta}}{d\varphi_\ell} = i\delta \sum_i G_{\ell,j,i}^{(1)} \phi_i(\mathbf{x}_\delta) - \frac{1}{2}\delta^2 \sum_{i,\ell'} \varphi_{\ell'} (G_{\ell\ell',ji}^{(2)} + G_{\ell'\ell,ji}^{(2)}) \phi_i(\mathbf{x}_\delta). \quad (158)$$

Moreover, we have

$$\begin{aligned} \dot{d}z_{j;\delta} &= i\delta \sum_i \sum_\ell \dot{d}\varphi_\ell G_{\ell,j,i}^{(1)} \phi_i(\mathbf{x}_\delta) - \frac{1}{2}\delta^2 \sum_i \sum_{\ell_1,\ell_2} \dot{d}\varphi_{\ell_1} \varphi_{\ell_2} G_{\ell_1\ell_2,ji}^{(2)} \phi_i(\mathbf{x}_\delta) \\ &\quad - \frac{1}{2}\delta^2 \sum_i \sum_{\ell_1,\ell_2} \varphi_{\ell_1} \dot{d}\varphi_{\ell_2} G_{\ell_1\ell_2,ji}^{(2)} \phi_i(\mathbf{x}_\delta) - \frac{1}{2}\delta^2 \sum_i \sum_{\ell_1,\ell_2} \dot{d}\varphi_{\ell_1} \dot{d}\varphi_{\ell_2} G_{\ell_1\ell_2,ji}^{(2)} \phi_i(\mathbf{x}_\delta). \end{aligned} \quad (159)$$

The leading order piece is exactly the effective feature map. A more compact version is,

$$\dot{d}z_{j;\delta} = \sum_\ell \frac{dz_{j;\delta}}{d\varphi_\ell} \dot{d}\varphi_\ell - \frac{1}{2}\delta^2 \sum_i \sum_{\ell_1,\ell_2} \dot{d}\varphi_{\ell_1} \dot{d}\varphi_{\ell_2} G_{\ell_1\ell_2,ji}^{(2)} \phi_i(\mathbf{x}_\delta). \quad (160)$$

Firstly we need to compute the effective kernel, we have

$$\begin{aligned} \frac{dz_{j;\delta}}{d\varphi_\ell} &= i\delta \sum_i G_{\ell,j,i}^{(1)} \phi_i(\mathbf{x}_\delta) - \frac{1}{2}\delta^2 \sum_{i,\ell'} \varphi_{\ell'} (G_{\ell\ell',ji}^{(2)} + G_{\ell'\ell,ji}^{(2)}) \phi_i(\mathbf{x}_\delta), \\ \frac{dz_{i';\bar{\alpha}}}{d\varphi_\ell} &= i\delta \sum_i G_{\ell,i',i}^{(1)} \phi_i(\mathbf{x}_{\bar{\alpha}}) - \frac{1}{2}\delta^2 \sum_{i,\ell'} \varphi_{\ell'} (G_{\ell\ell',i'i}^{(2)} + G_{\ell'\ell,i'i}^{(2)}) \phi_i(\mathbf{x}_{\bar{\alpha}}), \\ \frac{dz_{i';\bar{\alpha}}^*}{d\varphi_\ell} &= -i\delta \sum_i \phi_i^*(\mathbf{x}_{\bar{\alpha}}) G_{\ell,ii'}^{\dagger,(1)} - \frac{1}{2}\delta^2 \sum_{i,\ell'} \varphi_{\ell'} \phi_i^*(\mathbf{x}_{\bar{\alpha}}) (G_{\ell\ell',ii'}^{\dagger,(2)} + G_{\ell'\ell,ii'}^{\dagger,(2)}). \end{aligned} \quad (161)$$

Thus, we define

$$K_{\delta,\bar{\alpha}}^{E,\hat{\alpha}\hat{\beta},ii'} = \begin{pmatrix} K_{\delta,\bar{\alpha}}^{E,+,,ii'} & K_{\delta,\bar{\alpha}}^{E,-,,ii'} \\ K_{\delta,\bar{\alpha}}^{*,E,-,,ii'} & K_{\delta,\bar{\alpha}}^{*,E,+,,ii'} \end{pmatrix} (\varphi) = \begin{pmatrix} \sum_\ell \frac{dz_{i;\delta}}{d\varphi_\ell} \frac{dz_{i';\bar{\alpha}}^*}{d\varphi_\ell} \sum_\ell \frac{dz_{i;\delta}}{d\varphi_\ell} \frac{dz_{i';\bar{\alpha}}}{d\varphi_\ell} \\ \sum_\ell \frac{dz_{i;\delta}^*}{d\varphi_\ell} \frac{dz_{i';\bar{\alpha}}}{d\varphi_\ell} \sum_\ell \frac{dz_{i;\delta}}{d\varphi_\ell} \frac{dz_{i';\bar{\alpha}}}{d\varphi_\ell} \end{pmatrix}, \quad (162)$$

and

$$K_{\delta,\bar{\alpha}}^{E,\hat{\alpha}\hat{\beta},ii'}(0) = K_{\delta,\bar{\alpha}}^{\hat{\alpha}\hat{\beta},ii'} + K_{\delta,\bar{\alpha}}^{\Delta,\hat{\alpha}\hat{\beta},ii'} = \begin{pmatrix} K_{\delta,\bar{\alpha}}^{+,ii'} & K_{\delta,\bar{\alpha}}^{-,ii'} \\ K_{\delta,\bar{\alpha}}^{*,-,ii'} & K_{\delta,\bar{\alpha}}^{*,+,ii'} \end{pmatrix} + \begin{pmatrix} K_{\delta,\bar{\alpha}}^{\Delta,+,,ii'} & K_{\delta,\bar{\alpha}}^{\Delta,-,,ii'} \\ K_{\delta,\bar{\alpha}}^{*,\Delta,-,,ii'} & K_{\delta,\bar{\alpha}}^{*,\Delta,+,,ii'} \end{pmatrix}. \quad (163)$$

We find

$$\begin{aligned} K_{\delta,\bar{\alpha}}^{\Delta,+,,ii'} &= \frac{i}{2}\delta^3 \sum_{j,\ell',i''} \varphi_{\ell'}(0) (G_{\ell\ell',ij}^{(2)} + G_{\ell'\ell,ij}^{(2)}) G_{\ell,i''i'}^{\dagger,(1)} \phi_j(\mathbf{x}_\delta) \phi_{i''}^*(\mathbf{x}_{\bar{\alpha}}) \\ &\quad - \frac{i}{2}\delta^3 \sum_{j,\ell',i''} \varphi_{\ell'}(0) (G_{\ell\ell',i''i'}^{\dagger,(2)} + G_{\ell'\ell,i''i'}^{\dagger,(2)}) G_{\ell,ij}^{(1)} \phi_j(\mathbf{x}_\delta) \phi_{i''}^*(\mathbf{x}_{\bar{\alpha}}), \\ K_{\delta,\bar{\alpha}}^{\Delta,-,,ii'} &= -\frac{i}{2}\delta^3 \sum_{j,\ell',i''} \varphi_{\ell'}(0) (G_{\ell\ell',ij}^{(2)} + G_{\ell'\ell,ij}^{(2)}) G_{\ell,i''i''}^{(1)} \phi_j(\mathbf{x}_\delta) \phi_{i''}(\mathbf{x}_{\bar{\alpha}}) \\ &\quad - \frac{i}{2}\delta^3 \sum_{j,\ell',i''} \varphi_{\ell'}(0) (G_{\ell\ell',i''i''}^{(2)} + G_{\ell'\ell,i''i''}^{(2)}) G_{\ell,ij}^{(1)} \phi_j(\mathbf{x}_\delta) \phi_{i''}(\mathbf{x}_{\bar{\alpha}}). \end{aligned} \quad (164)$$

We could also write it in the bracket fashion

$$\begin{aligned} K_{\delta,\bar{\alpha}}^{\Delta,+,,ii'} &= \frac{i}{2}\delta^3 \sum_{\ell'} \varphi_{\ell'}(0) \begin{pmatrix} \langle i | (G_{\ell\ell'}^{(2)} + G_{\ell'\ell}^{(2)}) | \phi(\mathbf{x}_\delta) \rangle \langle \phi(\mathbf{x}_{\bar{\alpha}}) | G_{\ell}^{\dagger,(1)} | i' \rangle \\ - \langle \phi(\mathbf{x}_{\bar{\alpha}}) | (G_{\ell\ell'}^{\dagger,(2)} + G_{\ell'\ell}^{\dagger,(2)}) | i' \rangle \langle i | G_{\ell}^{(1)} | \phi(\mathbf{x}_\delta) \rangle \end{pmatrix}, \\ K_{\delta,\bar{\alpha}}^{\Delta,-,,ii'} &= -\frac{i}{2}\delta^3 \sum_{\ell'} \varphi_{\ell'}(0) \begin{pmatrix} \langle i | (G_{\ell\ell'}^{(2)} + G_{\ell'\ell}^{(2)}) | \phi(\mathbf{x}_\delta) \rangle \langle i' | G_{\ell}^{(1)} | \phi(\mathbf{x}_{\bar{\alpha}}) \rangle \\ + \langle i' | (G_{\ell\ell'}^{(2)} + G_{\ell'\ell}^{(2)}) | \phi_{i''}(\mathbf{x}_{\bar{\alpha}}) \rangle \langle i | G_{\ell}^{(1)} | \phi(\mathbf{x}_\delta) \rangle \end{pmatrix}. \end{aligned} \quad (165)$$

Now we could write down the prediction on the training set. We have

$$\varepsilon_{\hat{\mu}}(t) = \varepsilon_{\hat{\mu}}^F(t) + \varepsilon_{\hat{\mu}}^I(t), \quad (166)$$

where

$$\begin{aligned} \varepsilon_{\hat{\mu}}^F(t) &= \sum_{\hat{\mu}_1} U_{\hat{\mu}_1 \hat{\mu}_2}(t) \varepsilon_{\hat{\mu}_1}(0), \\ U_{\hat{\mu}_1 \hat{\mu}_2}(t) &= \left[ \left(1 - \frac{\eta}{2} K\right)^t \right]_{\hat{\mu}_1 \hat{\mu}_2}, \end{aligned} \quad (167)$$

and

$$\varepsilon_{\hat{\mu}}^I(t) = \left( -\frac{\eta}{2} \sum_{s=0}^{t-1} \left(1 - \frac{\eta}{2} K\right)^{t-1-s} K^\Delta \left(1 - \frac{\eta}{2} K\right)^s \varepsilon(0) \right)_{\hat{\mu}}. \quad (168)$$

It is the compact matrix product form in the space  $\mathcal{A} \times \mathbb{Z}_2 \times \mathcal{H}$ . Moreover, we have

$$\|\varepsilon^I(t)\| \leq \frac{\eta}{2} t \left\| 1 - \frac{\eta}{2} K \right\|^{t-1} \|K^\Delta\| \|\varepsilon(0)\|. \quad (169)$$

Finally, we could discuss the asymptotic convergence. We could start by computing the meta-kernel. We notice that

$$\begin{aligned} dz_{i;\delta} &= -\frac{\eta}{2} \sum_{\ell, i', \bar{\alpha}} \varepsilon_{i;\bar{\alpha}} \frac{dz_{i;\delta}}{d\varphi_\ell} \frac{dz_{i';\bar{\alpha}}^*}{d\varphi_\ell} - \frac{\eta}{2} \sum_{\ell, i', \bar{\alpha}} \varepsilon_{i';\bar{\alpha}}^* \frac{dz_{i;\delta}}{d\varphi_\ell} \frac{dz_{i';\bar{\alpha}}}{d\varphi_\ell} \\ &+ \left(\frac{\eta}{2}\right)^2 \sum_{\ell_1, \ell_2} \frac{d^2 z_{i;\delta}}{d\varphi_{\ell_1} d\varphi_{\ell_2}} \left( \sum_{\bar{\alpha}_1, i_1} \varepsilon_{i_1; \bar{\alpha}_1} \frac{dz_{i_1; \bar{\alpha}_1}^*}{d\varphi_{\ell_1}} + \sum_{\bar{\alpha}_1, i_1} \varepsilon_{i_1; \bar{\alpha}_1}^* \frac{dz_{i_1; \bar{\alpha}_1}}{d\varphi_{\ell_2}} \right) \left( \sum_{\bar{\alpha}_2, i_2} \varepsilon_{i_2; \bar{\alpha}_2} \frac{dz_{i_2; \bar{\alpha}_2}^*}{d\varphi_{\ell_2}} + \sum_{\bar{\alpha}_2, i_2} \varepsilon_{i_2; \bar{\alpha}_2}^* \frac{dz_{i_2; \bar{\alpha}_2}}{d\varphi_{\ell_2}} \right), \end{aligned} \quad (170)$$

We have

$$\begin{aligned} &\sum_{\ell_1, \ell_2} \frac{d^2 z_{i;\delta}}{d\varphi_{\ell_1} d\varphi_{\ell_2}} \left( \sum_{\bar{\alpha}_1, i_1} \varepsilon_{i_1; \bar{\alpha}_1} \frac{dz_{i_1; \bar{\alpha}_1}^*}{d\varphi_{\ell_1}} + \sum_{\bar{\alpha}_1, i_1} \varepsilon_{i_1; \bar{\alpha}_1}^* \frac{dz_{i_1; \bar{\alpha}_1}}{d\varphi_{\ell_2}} \right) \left( \sum_{\bar{\alpha}_2, i_2} \varepsilon_{i_2; \bar{\alpha}_2} \frac{dz_{i_2; \bar{\alpha}_2}^*}{d\varphi_{\ell_2}} + \sum_{\bar{\alpha}_2, i_2} \varepsilon_{i_2; \bar{\alpha}_2}^* \frac{dz_{i_2; \bar{\alpha}_2}}{d\varphi_{\ell_2}} \right) \\ &= \sum_{\ell_1, \ell_2, \bar{\alpha}_1, i_1, \bar{\alpha}_2, i_2, \dot{\alpha}_1, \dot{\alpha}_2} \frac{d^2 z_{i;\delta}}{d\varphi_{\ell_1} d\varphi_{\ell_2}} \frac{dz_{i_1; \bar{\alpha}_1}^{1-\dot{\alpha}_1}}{d\varphi_{\ell_1}} \frac{dz_{i_2; \bar{\alpha}_2}^{1-\dot{\alpha}_2}}{d\varphi_{\ell_2}} \varepsilon_{i_1; \bar{\alpha}_1}^{\dot{\alpha}_1} \varepsilon_{i_2; \bar{\alpha}_2}^{\dot{\alpha}_2}. \end{aligned} \quad (171)$$

So we define

$$\mu_{\delta \dot{\alpha}_1 \dot{\alpha}_2}^{ii_1 i_2; \dot{\beta} \dot{\alpha}_1 \dot{\alpha}_2} = \mu_{\bar{\mu} \hat{\mu}_1 \hat{\mu}_2} = \sum_{\ell_1, \ell_2} \frac{d^2 z_{i;\delta}^{\dot{\beta}}}{d\varphi_{\ell_1} d\varphi_{\ell_2}} \left( \frac{dz_{i_1; \bar{\alpha}_1}^{1-\dot{\alpha}_1}}{d\varphi_{\ell_1}} \frac{dz_{i_2; \bar{\alpha}_2}^{1-\dot{\alpha}_2}}{d\varphi_{\ell_2}} \right) \Big|_{\varphi=0}, \quad (172)$$

in the leading order. Moreover, in general we have,

$$\mu_{\delta_0 \dot{\delta}_1 \dot{\delta}_2}^{i_0 i_1 i_2; \dot{\alpha}_0 \dot{\alpha}_1 \dot{\alpha}_2} = \mu_{\bar{\mu}_0 \bar{\mu}_1 \bar{\mu}_2} = \sum_{\ell_1, \ell_2} \frac{d^2 z_{i_0; \delta_0}^{\dot{\alpha}_0}}{d\varphi_{\ell_1} d\varphi_{\ell_2}} \left( \frac{dz_{i_1; \bar{\delta}_1}^{1-\dot{\alpha}_1}}{d\varphi_{\ell_1}} \frac{dz_{i_2; \bar{\delta}_2}^{1-\dot{\alpha}_2}}{d\varphi_{\ell_2}} \right) \Big|_{\varphi=0}. \quad (173)$$

So the asymptotic convergence is given by

$$\begin{aligned} z_{\bar{\mu}}(\infty) &= z_{\bar{\mu}}(0) - \sum_{\hat{\mu}_1, \hat{\mu}_2} K_{\bar{\mu} \hat{\mu}_1} \tilde{K}^{\hat{\mu}_1 \hat{\mu}_2} \varepsilon_{\hat{\mu}_2}(0) \\ &+ \sum_{\hat{\mu}_1, \hat{\mu}_2, \hat{\mu}_3, \hat{\mu}_4} \left[ \mu_{\hat{\mu}_1 \bar{\mu} \hat{\mu}_2} - \sum_{\hat{\mu}_5, \hat{\mu}_6} K_{\bar{\mu} \hat{\mu}_5} \tilde{K}^{\hat{\mu}_5 \hat{\mu}_6} \mu_{\hat{\mu}_1 \hat{\mu}_6 \hat{\mu}_2} \right] Z_A^{\hat{\mu}_1 \hat{\mu}_2 \hat{\mu}_3 \hat{\mu}_4} \varepsilon_{\hat{\mu}_3}(0) \varepsilon_{\hat{\mu}_4}(0) \\ &+ \sum_{\hat{\mu}_1, \hat{\mu}_2, \hat{\mu}_3, \hat{\mu}_4} \left[ \mu_{\bar{\mu} \hat{\mu}_1 \hat{\mu}_2} - \sum_{\hat{\mu}_5, \hat{\mu}_6} K_{\bar{\mu} \hat{\mu}_5} \tilde{K}^{\hat{\mu}_5 \hat{\mu}_6} \mu_{\hat{\mu}_6 \hat{\mu}_1 \hat{\mu}_2} \right] Z_B^{\hat{\mu}_1 \hat{\mu}_2 \hat{\mu}_3 \hat{\mu}_4} \varepsilon_{\hat{\mu}_3}(0) \varepsilon_{\hat{\mu}_4}(0), \end{aligned} \quad (174)$$

where the algorithm projector  $Z_{A,B,S}$  are

$$\begin{aligned} Z_A^{\hat{\mu}_1 \hat{\mu}_2 \hat{\mu}_3 \hat{\mu}_4} &\equiv \tilde{K}^{\hat{\mu}_1 \hat{\mu}_3} \tilde{K}^{\hat{\mu}_2 \hat{\mu}_4} - \sum_{\hat{\mu}_5} \tilde{K}^{\hat{\mu}_2 \hat{\mu}_5} X_{\parallel}^{\hat{\mu}_1 \hat{\mu}_5 \hat{\mu}_3 \hat{\mu}_4}, \\ Z_B^{\hat{\mu}_1 \hat{\mu}_2 \hat{\mu}_3 \hat{\mu}_4} &\equiv \tilde{K}^{\hat{\mu}_1 \hat{\mu}_3} \tilde{K}^{\hat{\mu}_2 \hat{\mu}_4} - \sum_{\hat{\mu}_5} \tilde{K}^{\hat{\mu}_2 \hat{\mu}_5} X_{\parallel}^{\hat{\mu}_1 \hat{\mu}_5 \hat{\mu}_3 \hat{\mu}_4} + \frac{\eta}{4} X_{\parallel}^{\hat{\mu}_1 \hat{\mu}_2 \hat{\mu}_3 \hat{\mu}_4}, \end{aligned} \quad (175)$$

and

$$X_{\parallel}^{\hat{\mu}_1 \hat{\mu}_2 \hat{\mu}_3 \hat{\mu}_4} = \sum_{s=0}^{\infty} \left[ \left(1 - \frac{\eta}{2} K\right)^s \right]_{\hat{\mu}_1 \hat{\mu}_3} \left[ \left(1 - \frac{\eta}{2} K\right)^s \right]_{\hat{\mu}_2 \hat{\mu}_4}, \quad (176)$$

which is defined implicitly as

$$\delta_{\hat{\mu}_5}^{\hat{\mu}_1} \delta_{\hat{\mu}_6}^{\hat{\mu}_2} = \sum_{\hat{\mu}_3, \hat{\mu}_4} X_{\parallel}^{\hat{\mu}_1 \hat{\mu}_2 \hat{\mu}_3 \hat{\mu}_4} \left( \tilde{K}_{\hat{\mu}_3 \hat{\mu}_5} \delta_{\hat{\mu}_4 \hat{\mu}_6} + \delta_{\hat{\mu}_3 \hat{\mu}_5} \tilde{K}_{\hat{\mu}_4 \hat{\mu}_6} - \frac{\eta}{2} \tilde{K}_{\hat{\mu}_3 \hat{\mu}_5} \tilde{K}_{\hat{\mu}_4 \hat{\mu}_6} \right). \quad (177)$$

Now we write down all components of  $\mu$ . We have

$$\begin{aligned} \frac{d^2 z_{i_0; \delta_0}}{d\varphi_{\ell_1} d\varphi_{\ell_2}} &= -\frac{1}{2} \delta^2 \sum_i \left( G_{\ell_1 \ell_2, i_0 i}^{(2)} + G_{\ell_2 \ell_1, i_0 i}^{(2)} \right) \phi_i(\mathbf{x}_{\delta_0}), \\ \frac{d^2 z_{i_0; \delta_0}^*}{d\varphi_{\ell_1} d\varphi_{\ell_2}} &= -\frac{1}{2} \delta^2 \sum_i \phi_i^*(\mathbf{x}_{\delta_0}) \left( G_{\ell_1 \ell_2, i i_0}^{\dagger, (2)} + G_{\ell_2 \ell_1, i i_0}^{\dagger, (2)} \right), \\ \frac{dz_{i_1; \delta_1}^*}{d\varphi_{\ell_1}} &= -i\delta \sum_{i'} \phi_{i'}^*(\mathbf{x}_{\delta_1}) G_{\ell_1, i' i_1}^{\dagger, (1)}, & \frac{dz_{i_1; \delta_1}}{d\varphi_{\ell_1}} &= i\delta \sum_{i'} G_{\ell_1, i_1 i'}^{(1)} \phi_{i'}(\mathbf{x}_{\delta_1}), \\ \frac{dz_{i_2; \delta_2}^*}{d\varphi_{\ell_2}} &= -i\delta \sum_{i''} \phi_{i''}^*(\mathbf{x}_{\delta_2}) G_{\ell_2, i'' i_2}^{\dagger, (1)}, & \frac{dz_{i_2; \delta_2}}{d\varphi_{\ell_2}} &= i\delta \sum_{i''} G_{\ell_2, i_2 i''}^{(1)} \phi_{i''}(\mathbf{x}_{\delta_2}). \end{aligned} \quad (178)$$

So

$$\begin{aligned} \mu_{\delta_0 \delta_1 \delta_2}^{i_0 i_1 i_2; 000} &= \frac{1}{2} \delta^4 \sum_{i, i', i'', \ell_1, \ell_2} \left( G_{\ell_1 \ell_2, i_0 i}^{(2)} + G_{\ell_2 \ell_1, i_0 i}^{(2)} \right) G_{\ell_1, i' i_1}^{\dagger, (1)} G_{\ell_2, i'' i_2}^{\dagger, (2)} \phi_i(\mathbf{x}_{\delta_0}) \phi_{i'}^*(\mathbf{x}_{\delta_1}) \phi_{i''}^*(\mathbf{x}_{\delta_2}), \\ \mu_{\delta_0 \delta_1 \delta_2}^{i_0 i_1 i_2; 001} &= -\frac{1}{2} \delta^4 \sum_{i, i', i'', \ell_1, \ell_2} \left( G_{\ell_1 \ell_2, i_0 i}^{(2)} + G_{\ell_2 \ell_1, i_0 i}^{(2)} \right) G_{\ell_1, i' i_1}^{\dagger, (1)} G_{\ell_2, i_2 i''}^{(2)} \phi_i(\mathbf{x}_{\delta_0}) \phi_{i'}^*(\mathbf{x}_{\delta_1}) \phi_{i''}(\mathbf{x}_{\delta_2}), \\ \mu_{\delta_0 \delta_1 \delta_2}^{i_0 i_1 i_2; 010} &= -\frac{1}{2} \delta^4 \sum_{i, i', i'', \ell_1, \ell_2} \left( G_{\ell_1 \ell_2, i_0 i}^{(2)} + G_{\ell_2 \ell_1, i_0 i}^{(2)} \right) G_{\ell_1, i_1 i'}^{(1)} G_{\ell_2, i'' i_2}^{\dagger, (2)} \phi_i(\mathbf{x}_{\delta_0}) \phi_{i'}(\mathbf{x}_{\delta_1}) \phi_{i''}^*(\mathbf{x}_{\delta_2}), \\ \mu_{\delta_0 \delta_1 \delta_2}^{i_0 i_1 i_2; 011} &= \frac{1}{2} \delta^4 \sum_{i, i', i'', \ell_1, \ell_2} \left( G_{\ell_1 \ell_2, i_0 i}^{(2)} + G_{\ell_2 \ell_1, i_0 i}^{(2)} \right) G_{\ell_1, i_1 i'}^{(1)} G_{\ell_2, i_2 i''}^{(2)} \phi_i(\mathbf{x}_{\delta_0}) \phi_{i'}(\mathbf{x}_{\delta_1}) \phi_{i''}(\mathbf{x}_{\delta_2}), \\ \mu_{\delta_0 \delta_1 \delta_2}^{i_0 i_1 i_2; 100} &= \frac{1}{2} \delta^4 \sum_{i, i', i'', \ell_1, \ell_2} \left( G_{\ell_1 \ell_2, i i_0}^{\dagger, (2)} + G_{\ell_2 \ell_1, i i_0}^{\dagger, (2)} \right) G_{\ell_1, i' i_1}^{\dagger, (1)} G_{\ell_2, i'' i_2}^{\dagger, (2)} \phi_i^*(\mathbf{x}_{\delta_0}) \phi_{i'}^*(\mathbf{x}_{\delta_1}) \phi_{i''}^*(\mathbf{x}_{\delta_2}), \\ \mu_{\delta_0 \delta_1 \delta_2}^{i_0 i_1 i_2; 101} &= -\frac{1}{2} \delta^4 \sum_{i, i', i'', \ell_1, \ell_2} \left( G_{\ell_1 \ell_2, i i_0}^{\dagger, (2)} + G_{\ell_2 \ell_1, i i_0}^{\dagger, (2)} \right) G_{\ell_1, i' i_1}^{\dagger, (1)} G_{\ell_2, i_2 i''}^{(2)} \phi_i^*(\mathbf{x}_{\delta_0}) \phi_{i'}^*(\mathbf{x}_{\delta_1}) \phi_{i''}(\mathbf{x}_{\delta_2}), \\ \mu_{\delta_0 \delta_1 \delta_2}^{i_0 i_1 i_2; 110} &= -\frac{1}{2} \delta^4 \sum_{i, i', i'', \ell_1, \ell_2} \left( G_{\ell_1 \ell_2, i i_0}^{\dagger, (2)} + G_{\ell_2 \ell_1, i i_0}^{\dagger, (2)} \right) G_{\ell_1, i_1 i'}^{(1)} G_{\ell_2, i'' i_2}^{\dagger, (2)} \phi_i^*(\mathbf{x}_{\delta_0}) \phi_{i'}(\mathbf{x}_{\delta_1}) \phi_{i''}^*(\mathbf{x}_{\delta_2}), \\ \mu_{\delta_0 \delta_1 \delta_2}^{i_0 i_1 i_2; 111} &= \frac{1}{2} \delta^4 \sum_{i, i', i'', \ell_1, \ell_2} \left( G_{\ell_1 \ell_2, i i_0}^{\dagger, (2)} + G_{\ell_2 \ell_1, i i_0}^{\dagger, (2)} \right) G_{\ell_1, i_1 i'}^{(1)} G_{\ell_2, i_2 i''}^{(2)} \phi_i^*(\mathbf{x}_{\delta_0}) \phi_{i'}(\mathbf{x}_{\delta_1}) \phi_{i''}(\mathbf{x}_{\delta_2}). \end{aligned} \quad (179)$$

One could also write them in the bracket notations.

$$\begin{aligned}
\mu_{\delta_0 \delta_1 \delta_2}^{i_0 i_1 i_2; 000} &= \frac{1}{2} \delta^4 \sum_{\ell_1, \ell_2} \langle i_0 | \left( G_{\ell_1 \ell_2}^{(2)} + G_{\ell_2 \ell_1}^{(2)} \right) | \phi(\mathbf{x}_{\delta_0}) \rangle \langle \phi(\mathbf{x}_{\delta_1}) | G_{\ell_1}^{\dagger, (1)} | i_1 \rangle \langle \phi(\mathbf{x}_{\delta_2}) | G_{\ell_2}^{\dagger, (1)} | i_2 \rangle, \\
\mu_{\delta_0 \delta_1 \delta_2}^{i_0 i_1 i_2; 001} &= -\frac{1}{2} \delta^4 \sum_{\ell_1, \ell_2} \langle i_0 | \left( G_{\ell_1 \ell_2}^{(2)} + G_{\ell_2 \ell_1}^{(2)} \right) | \phi(\mathbf{x}_{\delta_0}) \rangle \langle \phi(\mathbf{x}_{\delta_1}) | G_{\ell_1}^{\dagger, (1)} | i_1 \rangle \langle i_2 | G_{\ell_2}^{(1)} | \phi(\mathbf{x}_{\delta_2}) \rangle, \\
\mu_{\delta_0 \delta_1 \delta_2}^{i_0 i_1 i_2; 010} &= -\frac{1}{2} \delta^4 \sum_{\ell_1, \ell_2} \langle i_0 | \left( G_{\ell_1 \ell_2}^{(2)} + G_{\ell_2 \ell_1}^{(2)} \right) | \phi(\mathbf{x}_{\delta_0}) \rangle \langle i_1 | G_{\ell_1}^{(1)} | \phi(\mathbf{x}_{\delta_1}) \rangle \langle \phi(\mathbf{x}_{\delta_2}) | G_{\ell_2}^{\dagger, (1)} | i_2 \rangle, \\
\mu_{\delta_0 \delta_1 \delta_2}^{i_0 i_1 i_2; 011} &= \frac{1}{2} \delta^4 \sum_{\ell_1, \ell_2} \langle i_0 | \left( G_{\ell_1 \ell_2}^{(2)} + G_{\ell_2 \ell_1}^{(2)} \right) | \phi(\mathbf{x}_{\delta_0}) \rangle \langle i_1 | G_{\ell_1}^{(1)} | \phi(\mathbf{x}_{\delta_1}) \rangle \langle i_2 | G_{\ell_2}^{(1)} | \phi(\mathbf{x}_{\delta_2}) \rangle, \\
\mu_{\delta_0 \delta_1 \delta_2}^{i_0 i_1 i_2; 100} &= \frac{1}{2} \delta^4 \sum_{\ell_1, \ell_2} \langle \phi(\mathbf{x}_{\delta_0}) | \left( G_{\ell_1 \ell_2}^{\dagger, (2)} + G_{\ell_2 \ell_1}^{\dagger, (2)} \right) | i_0 \rangle \langle \phi(\mathbf{x}_{\delta_1}) | G_{\ell_1}^{\dagger, (1)} | i_1 \rangle \langle \phi(\mathbf{x}_{\delta_2}) | G_{\ell_2}^{\dagger, (1)} | i_2 \rangle, \\
\mu_{\delta_0 \delta_1 \delta_2}^{i_0 i_1 i_2; 101} &= -\frac{1}{2} \delta^4 \sum_{\ell_1, \ell_2} \langle \phi(\mathbf{x}_{\delta_0}) | \left( G_{\ell_1 \ell_2}^{\dagger, (2)} + G_{\ell_2 \ell_1}^{\dagger, (2)} \right) | i_0 \rangle \langle \phi(\mathbf{x}_{\delta_1}) | G_{\ell_1}^{\dagger, (1)} | i_1 \rangle \langle i_2 | G_{\ell_2}^{(1)} | \phi(\mathbf{x}_{\delta_2}) \rangle, \\
\mu_{\delta_0 \delta_1 \delta_2}^{i_0 i_1 i_2; 110} &= -\frac{1}{2} \delta^4 \sum_{\ell_1, \ell_2} \langle \phi(\mathbf{x}_{\delta_0}) | \left( G_{\ell_1 \ell_2}^{\dagger, (2)} + G_{\ell_2 \ell_1}^{\dagger, (2)} \right) | i_0 \rangle \langle i_1 | G_{\ell_1}^{(1)} | \phi(\mathbf{x}_{\delta_1}) \rangle \langle \phi(\mathbf{x}_{\delta_2}) | G_{\ell_2}^{\dagger, (1)} | i_2 \rangle, \\
\mu_{\delta_0 \delta_1 \delta_2}^{i_0 i_1 i_2; 111} &= \frac{1}{2} \delta^4 \sum_{\ell_1, \ell_2} \langle \phi(\mathbf{x}_{\delta_0}) | \left( G_{\ell_1 \ell_2}^{\dagger, (2)} + G_{\ell_2 \ell_1}^{\dagger, (2)} \right) | i_0 \rangle \langle i_1 | G_{\ell_1}^{(1)} | \phi(\mathbf{x}_{\delta_1}) \rangle \langle i_2 | G_{\ell_2}^{(1)} | \phi(\mathbf{x}_{\delta_2}) \rangle. \tag{180}
\end{aligned}$$

#### D. Reading the amplitude

Here, we review the amplitude protocol for realizing the evaluation of the inner product of two states:  $\langle x|y \rangle$  [7]. The protocol is very similar to the celebrated Hadamard test. For a given pair of quantum states  $|x\rangle$  and  $|y\rangle$ , we need to get access to the state

$$|\varphi\rangle = \frac{1}{\sqrt{2}}(|0\rangle|x\rangle + |1\rangle|y\rangle). \tag{181}$$

The first qubit is serving as an ancillary qubit. Applying the Hadamard gate to the first qubit, we get

$$|\varphi\rangle \rightarrow \frac{1}{2}(|0\rangle(|x\rangle + |y\rangle) + |1\rangle(|x\rangle - |y\rangle)). \tag{182}$$

Now we could measure the probability to obtain  $|0\rangle$  in the ancillary system. We have

$$p = \frac{1}{2}(1 + \text{Re}(\langle x|y \rangle)). \tag{183}$$

Thus, we could use the probability to estimate the real part of the inner product. Moreover, we could add a phase rotation to  $|\varphi\rangle$  to get

$$|\varphi\rangle = \frac{1}{\sqrt{2}}(|0\rangle|x\rangle - i|1\rangle|y\rangle). \tag{184}$$

Then, we apply the same Hadamard gate and measure the first qubit in the Pauli- $Z$  basis. We get the probability,

$$p = \frac{1}{2}(1 + \text{Im}(\langle x|y \rangle)), \tag{185}$$

to obtain  $|0\rangle$ . Similar to the Hadamard test where we need to get access to a controlled unitary acting on an arbitrary state, here we need to get access to the state  $|\varphi\rangle$ . The inner product evaluation operation has a statistical error coming from the measurement. If we measure  $N$  times, the error scales as  $1/\sqrt{N}$ .

## V. SUPPRESSION OF NON-GAUSSIANITY IN THE LARGE-WIDTH LIMIT

In this section, we visit the statistics of hybrid quantum-classical neural networks. The model is defined as the following. First, we initialize the neural network by a quantum model,

$$z_{1;\alpha;j_1}^Q = \langle \phi_1(\mathbf{x}_\alpha) | U^{\dagger,1}(\theta^1) O_{j_1}^1 U^1(\theta^1) | \phi_1(\mathbf{x}_\alpha) \rangle. \quad (186)$$

Here we use the notation  $j_\omega$  to denote the index of the operator space  $\mathcal{O}^\omega(\mathcal{H}^\omega)$ , where  $\omega \in [1, 2, \dots, \Omega]$  is denoting the layer of the hybrid quantum-classical neural network. Here we are starting our first layer, so  $\omega = 1$ . We use  $z_{\omega;\alpha;j_\omega}^Q$  to denote our quantum model output. We also use

$$U^\omega(\theta^\omega) = \prod_{\ell_\omega=1}^{L_\omega} W_{\ell_\omega}^\omega \exp(i\theta_{\ell_\omega}^\omega X_{\ell_\omega}^\omega), \quad (187)$$

to denote our  $\omega$ th quantum ansatz, and  $\phi_\omega$  is used to denote the  $\omega$ th feature map. In each layer, after the quantum network, we connect it with a classical neural network by

$$w_{\omega;\alpha;j_\omega}^C = \sigma_{j_\omega}^C \left( \sum_{j_\omega=1}^{\dim \mathcal{O}^\omega(\mathcal{H}^\omega)} W_{j_\omega}^\omega z_{\omega;\alpha;j_\omega}^Q + b_{j_\omega}^C \right) \equiv \sigma_{j_\omega}^C(z_{\omega;\alpha;j_\omega}^C). \quad (188)$$

Here, we call

$$z_{\omega;\alpha;j_\omega}^C = \sum_{j_\omega=1}^{\dim \mathcal{O}^\omega(\mathcal{H}^\omega)} W_{j_\omega}^\omega z_{\omega;\alpha;j_\omega}^Q + b_{j_\omega}^C, \quad (189)$$

as classical preactivation in each layer, and we use the non-linear activation  $\sigma$ . At initialization, we will set all  $W$ s and  $b$ s distributed randomly from the following Gaussian statistics:

$$\begin{aligned} \mathbb{E} \left( W_{j_{1,\omega}^C, j_{1,\omega}}^\omega W_{j_{2,\omega}^C, j_{2,\omega}}^\omega \right) &= \delta_{j_{1,\omega}^C, j_{2,\omega}^C} \delta_{j_{1,\omega}, j_{2,\omega}} \frac{C_W^\omega}{\dim \mathcal{O}^\omega(\mathcal{H}^\omega)}, \\ \mathbb{E} \left( b_{j_{1,\omega}^C}^\omega b_{j_{2,\omega}^C}^\omega \right) &= \delta_{j_{1,\omega}^C, j_{2,\omega}^C} C_b^\omega. \end{aligned} \quad (190)$$

Moreover, after the classical network in each layer, we can move to the next quantum layer by doing the following encoding,

$$z_{\omega;\alpha;j_\omega}^Q = \langle \phi_\omega(\mathbf{w}_{\omega-1;\alpha}) | U^{\dagger,\omega}(\theta^\omega) O_{j_\omega}^\omega U^\omega(\theta^\omega) | \phi_\omega(\mathbf{w}_{\omega-1;\alpha}) \rangle, \quad (191)$$

where we use the vector notation  $\mathbf{w}_{\omega;\alpha} = (w_{\omega;\alpha})_{j_\omega}^C$ . Moreover, we will assume that the initial distribution of quantum variational angles is given by a statistical ensemble. We denote all those ensemble averages as  $\mathbb{E}$ . One could, for instance, compute the two-point function as

$$\begin{aligned} \mathbb{E} \left( z_{\omega;\alpha;j_{1,\omega}^C}^C z_{\omega;\beta;j_{2,\omega}^C}^C \right) &= \mathbb{E} \left( \left( \sum_{j_{1,\omega}=1}^{\dim \mathcal{O}^\omega(\mathcal{H}^\omega)} W_{j_{1,\omega}^C, j_{1,\omega}}^\omega z_{\omega;\alpha;j_{1,\omega}}^Q + b_{j_{1,\omega}^C}^\omega \right) \left( \sum_{j_{2,\omega}=1}^{\dim \mathcal{O}^\omega(\mathcal{H}^\omega)} W_{j_{2,\omega}^C, j_{2,\omega}}^\omega z_{\omega;\beta;j_{2,\omega}}^Q + b_{j_{2,\omega}^C}^\omega \right) \right) \\ &= \sum_{j_{1,\omega}, j_{2,\omega}=1}^{\dim \mathcal{O}^\omega(\mathcal{H}^\omega)} \mathbb{E} \left( W_{j_{1,\omega}^C, j_{1,\omega}}^\omega W_{j_{2,\omega}^C, j_{2,\omega}}^\omega \right) \mathbb{E} \left( z_{\omega;\alpha;j_{1,\omega}}^Q z_{\omega;\beta;j_{2,\omega}}^Q \right) + \mathbb{E} \left( b_{j_{1,\omega}^C}^\omega b_{j_{2,\omega}^C}^\omega \right) \\ &= \frac{C_W^\omega}{\dim \mathcal{O}^\omega(\mathcal{H}^\omega)} \sum_{j_{1,\omega}, j_{2,\omega}=1}^{\dim \mathcal{O}^\omega(\mathcal{H}^\omega)} \delta_{j_{1,\omega}^C, j_{2,\omega}^C} \delta_{j_{1,\omega}, j_{2,\omega}} \mathbb{E} \left( z_{\omega;\alpha;j_{1,\omega}}^Q z_{\omega;\beta;j_{2,\omega}}^Q \right) + \delta_{j_{1,\omega}^C, j_{2,\omega}^C} C_b^\omega \\ &= \delta_{j_{1,\omega}^C, j_{2,\omega}^C} \left( \frac{C_W^\omega}{\dim \mathcal{O}^\omega(\mathcal{H}^\omega)} \sum_{j_\omega=1}^{\dim \mathcal{O}^\omega(\mathcal{H}^\omega)} \mathbb{E} \left( z_{\omega;\alpha;j_\omega}^Q z_{\omega;\beta;j_\omega}^Q \right) + C_b^\omega \right). \end{aligned} \quad (192)$$

In order to compute non-Gaussianities of preactivations, we start from the connected part of the two-point function, which is given by

$$\begin{aligned} & \mathbb{E}_{\text{conn}} \left( z_{\omega; \alpha_1; j_{1,\omega}^C}^C z_{\omega; \alpha_2; j_{2,\omega}^C}^C z_{\omega; \alpha_3; j_{3,\omega}^C}^C z_{\omega; \alpha_4; j_{4,\omega}^C}^C \right) \\ & \equiv \mathbb{E} \left( z_{\omega; \alpha_1; j_{1,\omega}^C}^C z_{\omega; \alpha_2; j_{2,\omega}^C}^C z_{\omega; \alpha_3; j_{3,\omega}^C}^C z_{\omega; \alpha_4; j_{4,\omega}^C}^C \right) - \mathbb{E} \left( z_{\omega; \alpha_1; j_{1,\omega}^C}^C z_{\omega; \alpha_2; j_{2,\omega}^C}^C \right) \mathbb{E} \left( z_{\omega; \alpha_3; j_{3,\omega}^C}^C z_{\omega; \alpha_4; j_{4,\omega}^C}^C \right) \\ & - \mathbb{E} \left( z_{\omega; \alpha_1; j_{1,\omega}^C}^C z_{\omega; \alpha_3; j_{3,\omega}^C}^C \right) \mathbb{E} \left( z_{\omega; \alpha_2; j_{2,\omega}^C}^C z_{\omega; \alpha_4; j_{4,\omega}^C}^C \right) - \mathbb{E} \left( z_{\omega; \alpha_1; j_{1,\omega}^C}^C z_{\omega; \alpha_4; j_{4,\omega}^C}^C \right) \mathbb{E} \left( z_{\omega; \alpha_2; j_{2,\omega}^C}^C z_{\omega; \alpha_3; j_{3,\omega}^C}^C \right). \end{aligned} \quad (193)$$

We proceed by direct computation. We have

$$\begin{aligned} & \mathbb{E} \left( z_{\omega; \alpha_1; j_{1,\omega}^C}^C z_{\omega; \alpha_2; j_{2,\omega}^C}^C z_{\omega; \alpha_3; j_{3,\omega}^C}^C z_{\omega; \alpha_4; j_{4,\omega}^C}^C \right) \\ & = \mathbb{E} \left( \left( \sum_{j_{1,\omega}=1}^{\dim \mathcal{O}^\omega(\mathcal{H}^\omega)} W_{j_{1,\omega}^C, j_{1,\omega}}^\omega z_{\omega; \alpha_1; j_{1,\omega}}^Q + b_{j_{1,\omega}^C}^\omega \right) \left( \sum_{j_{2,\omega}=1}^{\dim \mathcal{O}^\omega(\mathcal{H}^\omega)} W_{j_{2,\omega}^C, j_{2,\omega}}^\omega z_{\omega; \alpha_2; j_{2,\omega}}^Q + b_{j_{2,\omega}^C}^\omega \right) \right) \\ & \left( \sum_{j_{3,\omega}=1}^{\dim \mathcal{O}^\omega(\mathcal{H}^\omega)} W_{j_{3,\omega}^C, j_{3,\omega}}^\omega z_{\omega; \alpha_3; j_{3,\omega}}^Q + b_{j_{3,\omega}^C}^\omega \right) \left( \sum_{j_{4,\omega}=1}^{\dim \mathcal{O}^\omega(\mathcal{H}^\omega)} W_{j_{4,\omega}^C, j_{4,\omega}}^\omega z_{\omega; \alpha_4; j_{4,\omega}}^Q + b_{j_{4,\omega}^C}^\omega \right) \\ & = \mathbb{E} \left( \sum_{j_{1,\omega}, j_{2,\omega}, j_{3,\omega}, j_{4,\omega}=1}^{\dim \mathcal{O}^\omega(\mathcal{H}^\omega)} W_{j_{1,\omega}^C, j_{1,\omega}}^\omega W_{j_{2,\omega}^C, j_{2,\omega}}^\omega W_{j_{3,\omega}^C, j_{3,\omega}}^\omega W_{j_{4,\omega}^C, j_{4,\omega}}^\omega z_{\omega; \alpha_1; j_{1,\omega}}^Q z_{\omega; \alpha_2; j_{2,\omega}}^Q z_{\omega; \alpha_3; j_{3,\omega}}^Q z_{\omega; \alpha_4; j_{4,\omega}}^Q \right) \\ & + \mathbb{E} \left( b_{j_{3,\omega}^C}^\omega b_{j_{4,\omega}^C}^\omega \sum_{j_{1,\omega}, j_{2,\omega}=1}^{\dim \mathcal{O}^\omega(\mathcal{H}^\omega)} W_{j_{1,\omega}^C, j_{1,\omega}}^\omega W_{j_{2,\omega}^C, j_{2,\omega}}^\omega z_{\omega; \alpha_1; j_{1,\omega}}^Q z_{\omega; \alpha_2; j_{2,\omega}}^Q \right) + (5 \text{ perms.}) \\ & + \mathbb{E} \left( b_{j_{1,\omega}^C}^\omega b_{j_{2,\omega}^C}^\omega b_{j_{3,\omega}^C}^\omega b_{j_{4,\omega}^C}^\omega \right). \end{aligned} \quad (194)$$

The notation 5 perms means that  $\binom{4}{2} - 1 = 5$  permutations of indices 1,2,3,4. Now, combining with the disconnected four-point function result,

$$\begin{aligned} & \mathbb{E} \left( z_{\omega; \alpha_1; j_{1,\omega}^C}^C z_{\omega; \alpha_2; j_{2,\omega}^C}^C \right) \mathbb{E} \left( z_{\omega; \alpha_3; j_{3,\omega}^C}^C z_{\omega; \alpha_4; j_{4,\omega}^C}^C \right) \\ & = \delta_{j_{1,\omega}^C, j_{2,\omega}^C} \delta_{j_{3,\omega}^C, j_{4,\omega}^C} \left( \left( \frac{C_W^\omega}{\dim \mathcal{O}^\omega(\mathcal{H}^\omega)} \sum_{j_{1,\omega}=1}^{\dim \mathcal{O}^\omega(\mathcal{H}^\omega)} \mathbb{E} \left( z_{\omega; \alpha_1; j_{1,\omega}}^Q z_{\omega; \alpha_2; j_{1,\omega}}^Q \right) + C_b^\omega \right) \right. \\ & \left. \left( \frac{C_W^\omega}{\dim \mathcal{O}^\omega(\mathcal{H}^\omega)} \sum_{j_{2,\omega}=1}^{\dim \mathcal{O}^\omega(\mathcal{H}^\omega)} \mathbb{E} \left( z_{\omega; \alpha_3; j_{2,\omega}}^Q z_{\omega; \alpha_4; j_{2,\omega}}^Q \right) + C_b^\omega \right) \right) \\ & = \delta_{j_{1,\omega}^C, j_{2,\omega}^C} \delta_{j_{3,\omega}^C, j_{4,\omega}^C} \left( \left( \frac{C_W^\omega}{\dim \mathcal{O}^\omega(\mathcal{H}^\omega)} \right)^2 \left( \sum_{j_{1,\omega}=1}^{\dim \mathcal{O}^\omega(\mathcal{H}^\omega)} \mathbb{E} \left( z_{\omega; \alpha_1; j_{1,\omega}}^Q z_{\omega; \alpha_2; j_{1,\omega}}^Q \right) \right) \left( \sum_{j_{2,\omega}=1}^{\dim \mathcal{O}^\omega(\mathcal{H}^\omega)} \mathbb{E} \left( z_{\omega; \alpha_3; j_{2,\omega}}^Q z_{\omega; \alpha_4; j_{2,\omega}}^Q \right) \right) \right. \\ & \left. + \frac{C_W^\omega C_b^\omega}{\dim \mathcal{O}^\omega(\mathcal{H}^\omega)} \left( \sum_{j_{1,\omega}=1}^{\dim \mathcal{O}^\omega(\mathcal{H}^\omega)} \mathbb{E} \left( z_{\omega; \alpha_1; j_{1,\omega}}^Q z_{\omega; \alpha_2; j_{1,\omega}}^Q \right) + \sum_{j_{2,\omega}=1}^{\dim \mathcal{O}^\omega(\mathcal{H}^\omega)} \mathbb{E} \left( z_{\omega; \alpha_3; j_{2,\omega}}^Q z_{\omega; \alpha_4; j_{2,\omega}}^Q \right) \right) \right. \\ & \left. + (C_b^\omega)^2 \right), \end{aligned} \quad (195)$$

and its two other  $t$  (14-23) and  $u$ -channel (13-24) permutations, we have

$$\begin{aligned} & \mathbb{E}_{\text{conn}} \left( z_{\omega; \alpha_1; j_{1,\omega}^C}^C z_{\omega; \alpha_2; j_{2,\omega}^C}^C z_{\omega; \alpha_3; j_{3,\omega}^C}^C z_{\omega; \alpha_4; j_{4,\omega}^C}^C \right) \\ & = \left( \frac{C_W^\omega}{\dim \mathcal{O}^\omega(\mathcal{H}^\omega)} \right)^2 \left( \begin{aligned} & \delta_{j_{1,\omega}^C, j_{2,\omega}^C} \delta_{j_{3,\omega}^C, j_{4,\omega}^C} \sum_{j_{1,\omega}, j_{2,\omega}=1}^{\dim \mathcal{O}^\omega(\mathcal{H}^\omega)} \mathbb{E}_{\text{conn}} \left( z_{\omega; \alpha_1; j_{1,\omega}}^Q z_{\omega; \alpha_2; j_{1,\omega}}^Q z_{\omega; \alpha_3; j_{2,\omega}}^Q z_{\omega; \alpha_4; j_{2,\omega}}^Q \right) \\ & + \delta_{j_{1,\omega}^C, j_{3,\omega}^C} \delta_{j_{2,\omega}^C, j_{4,\omega}^C} \sum_{j_{1,\omega}, j_{2,\omega}=1}^{\dim \mathcal{O}^\omega(\mathcal{H}^\omega)} \mathbb{E}_{\text{conn}} \left( z_{\omega; \alpha_1; j_{1,\omega}}^Q z_{\omega; \alpha_3; j_{1,\omega}}^Q z_{\omega; \alpha_2; j_{2,\omega}}^Q z_{\omega; \alpha_4; j_{2,\omega}}^Q \right) \\ & + \delta_{j_{1,\omega}^C, j_{4,\omega}^C} \delta_{j_{2,\omega}^C, j_{3,\omega}^C} \sum_{j_{1,\omega}, j_{2,\omega}=1}^{\dim \mathcal{O}^\omega(\mathcal{H}^\omega)} \mathbb{E}_{\text{conn}} \left( z_{\omega; \alpha_1; j_{1,\omega}}^Q z_{\omega; \alpha_4; j_{1,\omega}}^Q z_{\omega; \alpha_2; j_{2,\omega}}^Q z_{\omega; \alpha_3; j_{2,\omega}}^Q \right) \end{aligned} \right). \end{aligned} \quad (196)$$

Here, the connected piece made by quantum circuits is given by

$$\begin{aligned} & \mathbb{E}_{\text{conn}} \left( z_{\omega; \alpha_1; j_{1,\omega}}^Q z_{\omega; \alpha_2; j_{1,\omega}}^Q z_{\omega; \alpha_3; j_{2,\omega}}^Q z_{\omega; \alpha_4; j_{2,\omega}}^Q \right) \\ & \equiv \mathbb{E} \left( z_{\omega; \alpha_1; j_{1,\omega}}^Q z_{\omega; \alpha_2; j_{1,\omega}}^Q z_{\omega; \alpha_3; j_{2,\omega}}^Q z_{\omega; \alpha_4; j_{2,\omega}}^Q \right) \\ & - \mathbb{E} \left( z_{\omega; \alpha_1; j_{1,\omega}}^Q z_{\omega; \alpha_2; j_{1,\omega}}^Q \right) \mathbb{E} \left( z_{\omega; \alpha_3; j_{2,\omega}}^Q z_{\omega; \alpha_4; j_{2,\omega}}^Q \right). \end{aligned} \quad (197)$$

Now, we note that since we have  $\mathcal{O} \left( (\dim \mathcal{O}^\omega(\mathcal{H}^\omega))^2 \right)$  terms in the sum, the connected part will at most scale as  $\mathcal{O}(1)$  since the quantum outputs are made by normalized states whose norms are bounded by 1. However, if we assume

$$\mathbb{E}_{\text{conn}} \left( z_{\omega; \alpha_1; j_{1,\omega}}^Q z_{\omega; \alpha_2; j_{1,\omega}}^Q z_{\omega; \alpha_3; j_{2,\omega}}^Q z_{\omega; \alpha_4; j_{2,\omega}}^Q \right) = \mathcal{O}(1) \times \delta_{j_{1,\omega}, j_{2,\omega}}, \quad (198)$$

and their permutations for all  $\omega$ s, we have

$$\begin{aligned} & \mathbb{E}_{\text{conn}} \left( z_{\omega; \alpha_1; j_{1,\omega}}^Q z_{\omega; \alpha_2; j_{1,\omega}}^Q z_{\omega; \alpha_3; j_{2,\omega}}^Q z_{\omega; \alpha_4; j_{2,\omega}}^Q \right) \\ & = \left( \frac{C_W^\omega}{\dim \mathcal{O}^\omega(\mathcal{H}^\omega)} \right)^2 \left( \delta_{j_{1,\omega}, j_{2,\omega}}^C \delta_{j_{3,\omega}, j_{4,\omega}}^C + \delta_{j_{1,\omega}, j_{3,\omega}}^C \delta_{j_{2,\omega}, j_{4,\omega}}^C + \delta_{j_{1,\omega}, j_{4,\omega}}^C \delta_{j_{2,\omega}, j_{3,\omega}}^C \right) \times \sum_{j_{1,\omega}, j_{2,\omega}=1}^{\dim \mathcal{O}^\omega(\mathcal{H}^\omega)} \delta_{j_{1,\omega}, j_{2,\omega}} \times \mathcal{O}(1) \\ & = \left( \frac{C_W^\omega}{\dim \mathcal{O}^\omega(\mathcal{H}^\omega)} \right)^2 \times \mathcal{O}(\dim \mathcal{O}^\omega(\mathcal{H}^\omega)) \times \left( \delta_{j_{1,\omega}, j_{2,\omega}}^C \delta_{j_{3,\omega}, j_{4,\omega}}^C + \delta_{j_{1,\omega}, j_{3,\omega}}^C \delta_{j_{2,\omega}, j_{4,\omega}}^C + \delta_{j_{1,\omega}, j_{4,\omega}}^C \delta_{j_{2,\omega}, j_{3,\omega}}^C \right) \\ & = \frac{(C_W^\omega)^2}{\dim \mathcal{O}^\omega(\mathcal{H}^\omega)} \times \mathcal{O}(1) \times \left( \delta_{j_{1,\omega}, j_{2,\omega}}^C \delta_{j_{3,\omega}, j_{4,\omega}}^C + \delta_{j_{1,\omega}, j_{3,\omega}}^C \delta_{j_{2,\omega}, j_{4,\omega}}^C + \delta_{j_{1,\omega}, j_{4,\omega}}^C \delta_{j_{2,\omega}, j_{3,\omega}}^C \right). \end{aligned} \quad (199)$$

Thus, the connected part is suppressed by the large width  $1/\dim \mathcal{O}^\omega(\mathcal{H}^\omega)$ .

The orthogonal condition we impose for quantum neural networks might be satisfied by random assumptions of choices of operators  $O$  or random variational ansätze (say, averaging over the Pauli group or the 1-design). Similar arguments could be made for the dynamical NTK. It is beyond our scope in this work to discuss how to connect the suppressed non-Gaussian correlations randomized over initialization, and the suppression of dNTK (dQNTK) in dynamics in the large width. Thus, we hereby make more general arguments and leave the detailed work in the future. Consider the schematic formula of dQNTK in the hybrid network

$$\mu_{\delta_0 \delta_1 \delta_2}^{i_0 i_1 i_2} = \sum_{\ell_1, \ell_2} \frac{d^2 z_{i_0; \delta_0}}{d\theta_{\ell_1} d\theta_{\ell_2}} \left( \frac{dz_{i_1; \delta_1}}{d\theta_{\ell_1}} \frac{dz_{i_2; \delta_2}}{d\theta_{\ell_2}} \right), \quad (200)$$

We schematically use notation  $\theta_\ell$  to denote both classical and quantum training variables. The sum will lead to a combination of the following three cases: classical contribution, quantum contribution, and the classical-quantum mixed contribution. For pure classical contribution, the formula from Gaussian correlations will lead to  $1/\text{width}$  suppression. For quantum contribution and its mixture with classical ones, the meta-kernel will be naturally suppressed by  $1/\dim \mathcal{O}(\mathcal{H})$  for its corresponding operator space dimension  $\dim \mathcal{O}(\mathcal{H})$ . Otherwise, there are non-negligible terms that give  $\mathcal{O}(1)$  modification to NTK during the gradient descent, leading to a non-linear regime of training dynamics and representation learning. Otherwise, the training dynamics will linearize, and we get an exponential convergence of the residual training error. The observations indicate a possible connection between barren plateau, random unitary, and large-width limit in classical neural networks, where we will leave details to our future works.

## VI. ON THE LECUN PARAMETRIZATION AND THE NTK PARAMETRIZATION

In this section, we will clarify the issue of the initialization (parametrization) convention we use comparing to different literature. In several original papers about classical NTK [8], there is a convention which is called the NTK parametrization, which defines the Gaussian correlation of weights as

$$\mathbb{E}(W_{ij} W_{kl}) = \delta_{ik} \delta_{jl} C_W. \quad (201)$$

The initialization will lead to rescaling of parameters as the following table,

In our work, we use the standard (LeCun) scaling, where we will still naturally obtain the natural NTK, and the dNTK suppression results [1]. However, there are also studies pointing out alternative parametrizations [9].

Type	Weight	Single-layer	Learning rate
LeCun	$\mathbb{E}(W_{ij}W_{kl}) = \delta_{ik}\delta_{jl}C_W/\text{width}$	$Wx + b$	$\eta = \mathcal{O}(1)/\text{width}$
NTK	$\mathbb{E}(W_{ij}W_{kl}) = \delta_{ik}\delta_{jl}C_W$	$(Wx + b)/\sqrt{\text{width}}$	$\eta = \mathcal{O}(1)$

TABLE I. Comparison between the LeCun (standard) parametrization and the NTK parametrization.

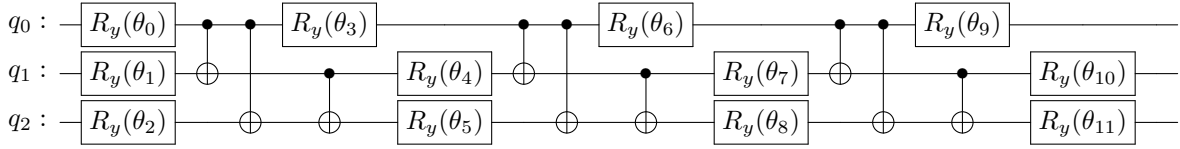
## VII. SIMULATION DETAILS AND ADDITIONAL DATA

We simulate the NTK dynamics in a concrete example using the Qiskit python library. In the main text, we consider a binary classification task carried out with a 3-qubit quantum neural network (QNN) in a supervised learning setting, using an ad-hoc data set provided in Qiskit Machine Learning. The input data points have 3-dimensional features, and are partitioned in a training set (20 elements) and a test set (5 elements). Each input  $\mathbf{x} = (x[0], x[1], x[2])$  is encoded through the ZZFeatureMap

$$\text{ZZFeatureMap} = \left. \begin{array}{l} q_0 : \left[ \text{H} \rightarrow \phi(2x_0) \right] \\ q_1 : \left[ \text{H} \rightarrow \phi(2x_1) \right] \\ q_2 : \left[ \text{H} \rightarrow \phi(2x_2) \right] \end{array} \right\}^2, \quad (202)$$

where  $\phi(\theta)$  is the single-qubit phase gate and  $\hat{x}_i = \pi - x_i$ . The trainable part of the quantum neural network, which appended to the quantum circuit after the feature map, is represented by a 3-layer RealAmplitudes variational ansatz, combining parametrized  $R_y$  rotations and entangling CNOT operations:

$$\text{RealAmplitudes} = \quad . \quad (203)$$



The QNN output is obtained as the expectation value of the 3-qubit observable  $Z_0Z_1Z_2$ , where  $Z_i$  is the single qubit Pauli  $Z$  operator. Our model is

$$z_\delta = \langle \phi(\mathbf{x}_\delta) | U^\dagger O U | \phi(\mathbf{x}_\delta) \rangle, \quad (204)$$

where  $O$  is given by

$$O = \text{diag}(1, -1, -1, 1, -1, 1, 1, -1). \quad (205)$$

Note that we only have one observable  $O$ , so our quantum neural network only has a scalar as the output. In this setup, the QNTK is given by

$$K_{\delta, \bar{\alpha}} = \sum_\ell \frac{dz_\delta}{d\theta_\ell} \frac{dz_{\bar{\alpha}}}{d\theta_\ell} = - \sum_\ell \left( \left\langle \phi(\mathbf{x}_\delta) \left| U_{+, \ell}^\dagger \left[ X_\ell, U_\ell^\dagger W_\ell^\dagger U_{-, \ell}^\dagger O U_{-, \ell} W_\ell U_\ell \right] U_{+, \ell} \right| \phi(\mathbf{x}_\delta) \right\rangle \times \left\langle \phi(\mathbf{x}_{\bar{\alpha}}) \left| U_{+, \ell}^\dagger \left[ X_\ell, U_\ell^\dagger W_\ell^\dagger U_{-, \ell}^\dagger O U_{-, \ell} W_\ell U_\ell \right] U_{+, \ell} \right| \phi(\mathbf{x}_{\bar{\alpha}}) \right\rangle \right), \quad (206)$$

and the frozen QNTK is given by

$$K_{\delta, \bar{\alpha}} = -\delta^2 \sum_\ell \left( \left\langle \phi(\mathbf{x}_\delta) \left| W_{+, \ell}^\dagger \left[ X_\ell, W_\ell^\dagger W_{-, \ell}^\dagger O W_{-, \ell} W_\ell \right] W_{+, \ell} \right| \phi(\mathbf{x}_\delta) \right\rangle \times \left\langle \phi(\mathbf{x}_{\bar{\alpha}}) \left| W_{+, \ell}^\dagger \left[ X_\ell, W_\ell^\dagger W_{-, \ell}^\dagger O W_{-, \ell} W_\ell \right] W_{+, \ell} \right| \phi(\mathbf{x}_{\bar{\alpha}}) \right\rangle \right), \quad (207)$$

with the help of the initial angle redefinition. All the theories we have discussed could be applied in a straightforward way in our simulation case. For the theoretical prediction, we take the variational angle  $\theta^*$  at the last step of our training.

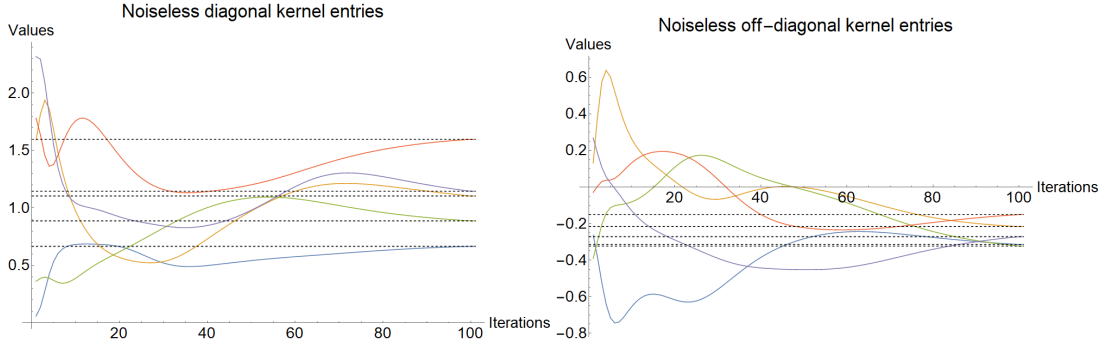


FIG. 1. Noiseless diagonal and off-diagonal entries of QNTK during the noiseless gradient descent dynamics. We compute the gradient descent evolution of five random diagonal and off-diagonal elements in the QNTK. The solid line is the actual value of the QNTK entries during the experiment, and the dashed line is the theoretical prediction of the frozen QNTK.

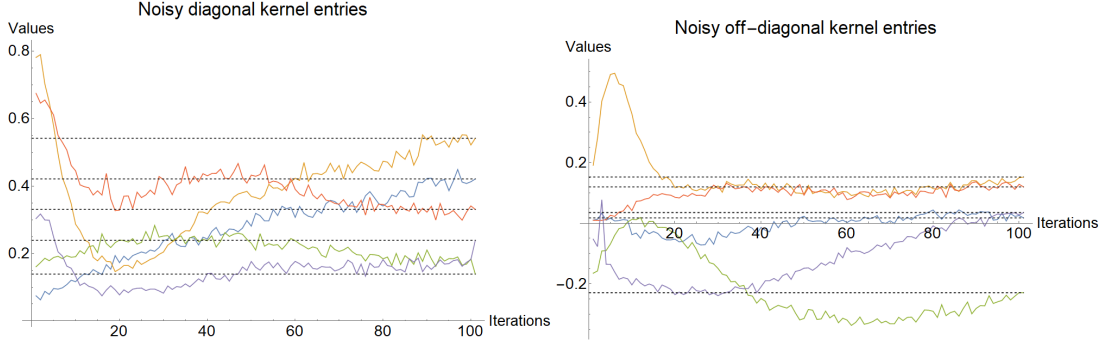


FIG. 2. Noisy diagonal and off-diagonal entries of QNTK during the gradient descent dynamics. We compute the gradient descent evolution of five random diagonal and off-diagonal elements in the QNTK. The solid line is the actual value of the QNTK entries during the simulation of a noisy quantum processor, and the dashed line is the theoretical prediction of the frozen QNTK.

In addition to the plots reported in the main text, here we show more data about the properties of QNTK during the gradient descent. See Fig. 1 and Fig. 2. We could see that, even including a model of device noise, the simulation is successful and the kernel is stable at the late time including the effect of the error mitigation. This will indicate an exponential decay of the residual training error at the late time for noisy quantum circuits.

- 
- [1] D. A. Roberts, S. Yaida, and B. Hanin, arXiv preprint arXiv:2106.10165 (2021).
  - [2] D. A. Roberts and S. Yaida, Deep Learning Theory Summer School at Princeton (2021).
  - [3] Note that the inverse has to be taken in the over-parametrized regime where the number of variational angles are larger than the training samples to make the formalism simple, which is not the case in our numerical simulation in this paper.
  - [4] D. Meltzer, H. Zheng, D. Yi-Hsien, D. R. Roberts, and J. Liu, .
  - [5] M. Mohri, A. Rostamizadeh, and A. Talwalkar, *Foundations of machine learning* (MIT press, 2018).
  - [6] The convention factor  $1/2$  appearing in  $\eta$  is from the worldsheet index defined for holomorphic and anti-holomorphic variables. Similar conventions appear in the two-dimensional conformal field theories, see [10].
  - [7] L. Zhao, Z. Zhao, P. Rebentrost, and J. Fitzsimons, Quantum Machine Intelligence **3**, 1 (2021).
  - [8] A. Jacot, F. Gabriel, and C. Hongler, arXiv preprint arXiv:1806.07572 (2018).
  - [9] J. Sohl-Dickstein, R. Novak, S. S. Schoenholz, and J. Lee, arXiv preprint arXiv:2001.07301 (2020).
  - [10] J. Polchinski, *String theory: Volume 1, an introduction to the bosonic string* (Cambridge university press, 1998).