Gauge covariant neural network for 4 dimensional non-abelian gauge theory

Akio Tomiya,1* Yuki Nagai^{2,3}

¹RIKEN/BNL Research center, Brookhaven National Laboratory, Upton, NY, 11973, USA,
 ²CCSE, Japan Atomic Energy Agency, 178-4-4, Wakashiba, Kashiwa, Chiba, 277-0871, Japan,
 ³Mathematical Science Team, RIKEN Center for Advanced Intelligence Project (AIP)
 1-4-1 Nihonbashi, Chuo-ku, Tokyo 103-0027, Japan

*E-mail: akio.tomiya@riken.jp

We develop a gauge covariant neural network for four dimensional non-abelian gauge theory, which realizes a map between rank-2 tensor valued vector fields. We find that the conventional smearing procedure and gradient flow for gauge fields can be regarded as known neural networks, residual networks and neural ordinal differential equations for rank-2 tensors with fixed parameters. In terms of machine learning context, projection or normalization functions in the smearing schemes correspond to an activation function in neural networks. Using the locality of the activation function, we derive the backpropagation for the gauge covariant neural network. Consequently, the smeared force in hybrid Monte Carlo (HMC) is naturally derived with the backpropagation. As a demonstration, we develop the self-learning HMC (SLHMC) with covariant neural network approximated action for non-abelian gauge theory with dynamical fermions, and we observe SLHMC reproduces results from HMC.

1 Introduction

Neural networks have been used among theoretical physics, condensed matter, and high energy physics (1, 2). Neural networks enable us to approximate functions in a systematic way (3-7), like variational ansatz of quantum states, and detector of phase boundaries (8-12).

A key concept in theoretical physics is symmetry. Historically, the most famous example of the power of symmetry is the derivation of special relativity. Assuming Lorentz symmetry in electromagnetism, which is a spacetime global symmetry, on mechanics, one can reach special relativity. In general relativity, spacetime local symmetries play essential roles. For quantum systems, long range behavior is determined by a global symmetry for fields. For example, continuous global symmetry leads Noether currents and conserved charges. If the symmetry is spontaneously broken, we expect the emergence of gapless modes. In condensed matter, symmetry protects a topology for wave functions in solids (13). Gauge symmetry, transformation depending on a coordinate, determines the dynamics of the theory. For example, gauge symmetry is quite important in the standard model of particle physics. The gauge symmetry is symbolically written $SU(3) \times SU(2) \times U(1)$ and it is the most successful model in physics (14).

In machine learning, symmetry is also important because of one of the domain knowledge. The performance of neural networks is improved if a parameter space is restricted by symmetry. Convolutional neural networks well deal with translation invariance, and it actually has group equivariant structure (15). This idea is generalized as spherical convolution (16). Spherical convolution helps to treat data with rotational symmetry in three dimension. Taco Cohen *et al.* invented treatment of a discrete gauge symmetry as a generalization of convolution (17) in the context of machine learning and computer vision.

Gauge theory is based on local symmetry, and naive discretization for numerical calculations breaks this symmetry. The lattice gauge theory (30) is one of the most powerful and quantita-

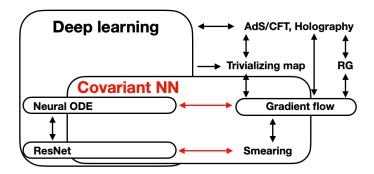


Figure 1: Black arrows indicate relations discussed in (17–29). See main text in details. RG, NN, and neural ODE indicate renormalization group, neural network, and neural ordinal differential equation, respectively. Red arrows and red text will be discussed in the present paper.

tive tools to investigate particle and nuclear physics. By defining gauge link (Wilson line) on bonds on discretized spacetime, one can obtain the lattice gauge theory. Lattice gauge theory enables us to obtain non-perturbative information of Quantum Chromo-Dynamics (QCD) through numerical calculations on supercomputers (31,32). Great achievements have done with a Markov chain Monte-Carlo algorithm called Hybrid Monte Carlo (HMC) (33) and its variant (34). We can calculate matrix elements in matrix elements for the anomalous magnetic moments in g-2 (35), form factors and low-energy constants for analytic calculation (36), and thermodynamics of QCD (37) using lattice QCD and its numerical calculations.

In the last twenty years, a lot of theoretical and numerical progress appeared in lattice QCD. M. Luscher introduced the gradient flow to lattice gauge theory away from differential geometry, which is a fictitious euclidean time evolution based on a gradient of a field theoretic action intended to build a trivializing map of gauge theory (18, 38). Later, he pointed out the gradient flow can be regarded as a continuous version of stout smearing (39). Gradient flow is not only useful for lattice gauge theory but also general field theory. For example, applying the gradient flow to the non-linear sigma model, one can obtain AdS geometry (26–29). Another example is relation to the renormalization group. Relations of the gradient flow and the renormalization

group have been discussed (40-43). Recently, Sonoda *et al.* found that flowed field can be regarded as smoothly regulated field respecting gauge symmetry, and one can define Polchinskitype renormalization group equation with keeping gauge invariance (42, 43).

In the last five years, the progress of machine learning has led us to apply it to lattice field theory. For example, finding new order parameters (11, 44-46), reduction of cost in valence calculations (47-49), and producing field configuration on the lattice (50-58).

Recently, the ideas of the fields of the lattice QCD and machine learning are being combined. Construction of the trivializing map using the gradient flow which is originated from Luscher's work has been continued with neural network (25, 54, 56-58). Trivializing map for lattice QCD has been demanded because once we obtain a trivializing map, we do not have to care autocorrelation in the simulation (18). However, the construction is not trivial and they tried to build it using neural networks. If we regard a neural network as a holographic renormalization flow, one can construct a trivializing map (25). On the other hand, the construction of a trivializing map has a problem in practice with the Jacobian, which requires $O(V^3)$ numerical cost. However, if one embeds symplectic type transformation in a neural network, calculation of determinant magically reduced (54, 56-58).

Physicists who were interested in deep learning tried to understand why deep learning framework from a physics point of view (19, 20, 59, 60). In particular, the similarity between deep neural networks and the renormalization group flow has been discussed. On the one hand, K. Hashimoto *et al.* regarded a deep neural network as an emergent spacetime (21–23, 61), especially a spacetime with AdS geometry. Moreover, a series of these work provides holographic renormalization picture (61).

In this paper, we find a novel connection between the machine learning technique and the lattice gauge theory. One of the most famous deep learning architecture is Res-Net (62, 63), which has a skip connection to preserve information during the forward propagation. This

particular structure enables us to regard it as a difference equation and consider a continuum limit of it. This idea is called neural-ordinal differential equation (neural ODE) (24). This is just a change of viewpoint, but it gives us an efficient way to perform backpropagation in training using the adjoint method. As we will explain later, the gradient flow can be regarded as a neural ODE of the covariant neural network with fixed parameters.

We develop gauge covariant neural network and gauge invariant loss function for lattice gauge theory as a component of machine learning technique for lattice QCD. We find that it is a generalization of smearing with tunable parameters *i.e.* smearing can be regarded as the gauge covariant neural network for forward propagation. There already similar relation there in the convolutional neural network, which can be regarded as a filter on an image with trainable parameters (Fig. 2). In addition, the delta rule in backpropagation for covariant neural networks leads HMC force formula in a systematic way. As a use case, we perform a simulation with self-learning hybrid Monte-Carlo (SLHMC) (64) for non-abelian gauge theory with dynamical fermion. This is a natural extension of our previous work (53) with self-learning Monte-Carlo (SLMC). SLHMC uses parametrized action in the molecular dynamics and we use output from gauge covariant neural network as an ansatz of the action. In contrast to SLMC for QCD (64), we do not have to evaluate fermion determinant and we simplify the Metropolis step because the molecular dynamics is time reversible.

The covariant neural networks can be regarded as a generalization of conventional neural networks. A conventional neural network is fed a vector and output is a vector. From a field theoretic point of view, what fed is, a set of scalar values, namely scalar field. In our framework, our network is fed matrices (rank-2 tensor valued vector fields), and they are transformed by repetitions of linear transformations and non-linear functions. As a output, we get a rank 2 tensor field, matrix field, which is consisted by a member of gauge group. This seems too abstract to construct but this is what happens in the smearing procedure and the Luscher's

gradient flow.

Figure 2: Convolutional neural network can be regarded as tunble filters for images. Analogously, covariant neural network can be regarded as smearing with tunble parameters for configurations.

There already similar works (17, 56–58, 65), namely neural network for gauge symmetric data because treatment of gauge symmetry was a problem (50). Early work by T. Cohen *et al.* (17) invented as an extension of group equivariant neural network (15, 66). MIT and Google brain group employ and extended these ideas on the flow based sampling method for gauge theory (56–58). Recently, to simulate many-body quantum systems efficiently with exact local gauge invariance, gauge equivariant neural-network quantum states are introduced (65). In stead of gauge equivariance, we implement gauge covariance¹ into a neural network. Our formulation is applicable for SU(N) in more than two dimensions as smearing.

The rest of the present paper is organized as follows. In the next section, we introduce a neural network with covariant layers, which is information processing of rank-2 vector fields. We discuss relations to well-known frameworks in lattice QCD and machine learning. To train the network, we derive the delta rule for the rank-2 neural network, which coincides with derivation for HMC force. We also discuss a connection to neural ODE. In the next section, we show results for a use case of the covariant neural networks in SLHMC. As a result, we get consistent results from HMC. This means that the covariant neural network succeeded to parametrize an action. In the last section, we summarize this work.

¹Concept of covariance is explained in Appendix A.2. In lattice gauge theory, a matter field $\psi(n)$ is transforms like $G(n)\psi(n)$, G(n) is a transform function of the gauge group, and a composite operator $\bar{\psi}(n)U_{\mu}(n)\psi(n+\hat{\mu})$ is invariant since $U_{\mu}(n)$ transforms covariantly. Please refer for gauge equivariance (67, 68).

2 Covariant Neural network

2.1 Overview and structure

2.1.1 Gauge covariant neural networks

Here we introduce covariant neural networks for lattice gauge theory². A covariant neural network realizes a map between gauge fields,

$$U \mapsto U_{\theta}^{\text{NN}}[U].$$
 (1)

with stacking of covariant layers. A covariant layer coincide with a smeared link at site n with

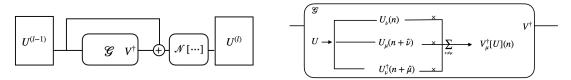


Figure 3: (Left) Calculation graph notation of neural networks for generalized smearing procedure. Input and output of this component are matrices. For stout type, \mathcal{N} is taken to be identity (see Fig. 4). (Right) Non-linear transformation \mathcal{G} for APE type is displayed, which constructs staples and Wilson lines which has same start and end point of $U_{\mu}(n)$. For stout-type, it is constructed by exponential function of un-traced plaquette.

direction μ with tunable parameters,

$$z_{\mu}^{(l)}(n) = w_1^{(l-1)} U_{\mu}^{(l-1)}(n) + w_2^{(l-1)} \mathcal{G}_{\mu,n}^{\bar{\theta}^{(l-1)}}(U^{(l-1)}), \tag{2}$$

$$U_{\mu}^{(l)}(n) = \mathcal{N}[z_{\mu}^{(l)}(n)],\tag{3}$$

where $\mathcal{N}(z)$ and $\mathcal{G}_{\mu,n}^{\bar{\theta}^{(l-1)}}$ are matrix-valued functions and $U_{\mu}^{(0)}(n) = U_{\mu}(n)$, and $z_{\mu}^{(l=0)}(n) = U_{\mu}(n)$. $\mathcal{G}_{\mu,n}^{\bar{\theta}^{(l-1)}}(U^{(l-1)})$ is sum of staples or Wilson lines which has same origin and end with a link $U_{\mu}^{(l-1)}(n)$ to keep gauge covariance (right panel of Fig. 3 for example) for APE-type case. In general, it can be a gauge covariant combination of links and $\bar{\theta}^{(l-1)}$ is a set of parameters in it. In conventional smearing scheme, as we will see later, \mathcal{N} is projection

²Conventional neural networks and their training are explain in Appendix A.1.

or normalization to $SU(N_c)$, which corresponds to an activation function because it acts on locally³. $w_1^{(l-1)}, w_2^{(l-1)}$ are real parameters, which have to be determine in training. In general, they can be taken as vector fields, $w_i^{(l-1)} \to w_{i,\mu}^{(l-1)}(n)$ and the current formulation corresponds to weight shared one⁴. As conventional smearing, we can stack this procedure obviously. From this point of view, multi steps of smearing corresponds to a deep neural network with fixed parameters and we will use this, to describe a relation to neural ODE. Symbolically, forward propagation of information is written as,

$$U = U^{(l=0)} \to U^{(l=1)}(w^{(l=0)}; U^{(0)}) \to U^{(l=2)}(w^{(l=1)}; U^{(1)}) \to \cdots$$
$$\to U^{(L)}(w^{(l=L-1)}; U^{L-1}) \equiv U_{\theta}^{\text{NN}}[U], \tag{4}$$

where we denote U_{θ}^{NN} , θ is a set of parameters, namely it is collection of $w_i^{(l)}$ and $\bar{\theta}^{(l)}$ (i=1,2 and $l=0,1,2,\cdots,L-1$). Or equivalently, to see connection to the conventional neural network Eq. (47), we can write it as,

$$U_{\theta}^{\text{NN}}[U] = U^{(L)}(U^{(L-1)}(U^{L-2}\cdots (U^{(1)}(U))\cdots)). \tag{5}$$

This is the covariant neural network. This transforms under a gauge transformation,

$$(U_{\theta}^{\text{NN}}[U])_{\mu}(n) \to G(n)(U_{\theta}^{\text{NN}}[U])_{\mu}(n)G(n+\hat{\mu}),\tag{6}$$

where $G(n + \hat{\mu})$ is a transformation function at site n, as same as the original link as in the conventinal smearing by construction.

Note that, in our definition, weight θ in the covariant neural network does not depend on on the volume because of the weight sharing feature. Thus we can expect that once trained in a small lattice, we are able to apply it to a larger lattice and adjustable with small number of training.

³This activation function \mathcal{N} has different locality to $\mathcal{G}_{\mu,n}^{\bar{\theta}}$. First one acts locally and this property corresponds to element-wise non-linear operation in conventional neural networks.

⁴From a viewpoint of geometrical symmetry, we should use weight sharing, which explicitly guarantees rotational and transnational symmetries.

After training, numerical cost of the covariant neural network is exactly same as the conventional smearing schemes.

2.1.2 Gauge invariant output and loss function

One of the usages of the covariant neural networks is for constructing effective action S_{θ} . The gauge invariant effective action S_{θ} is constructed by

$$S_{\theta} = \sum_{\mu,n} S_{\mu,n}[U_{\theta}^{\text{NN}}[U]] \tag{7}$$

where $S_{\mu,n}[U]$ is constructed by traced loop operators with smeared links $U_{\theta}^{\text{NN}}[U]$. If we regard $S_{\mu,n}[U]$ as a output of a neural network with input U, this network is the Behler-Parrinello type neural works which has been introduced in the machine-learning molecular dynamics (69, 70). While the Behler-Parrinello type neural work in condensed matter physics community uses a weight sharing between same kind of atomic species, this gauge covariant neural network uses a weight sharing between gauge fields on space time lattice.

Naturally, we can construct gauge invariant loss function using a gauge invariant field theoretic action. Namely

$$L_{\theta}(\mathcal{D}) = f(S_{\theta}, \mathcal{D}), \tag{8}$$

where $S_{\theta} = S(U_{\theta}^{\text{NN}}[U])$ and f is a conventional loss function like square of difference. \mathcal{D} is a data for training and we intend that value of action in the target system. S is a gauge invariant action, which is typically constructed by traced loop operators with smeared links $U_{\theta}^{\text{NN}}[U]$ and, we can use a fermionic action as in the conventional smearing.

As we will show later, this network can be trained using the backpropagation. In addition, HMC force can be evaluated by the backpropagation, it is indeed almost the same derivation for usual smearing.

2.2 Connection to smearing

In this subsection, we identify relation between the covariant neural network and smearing schemes.

2.2.1 Connection to APE-type smearing

For APE-type smearing (71, 72), the parameters and functions are taken as,

$$w_1^{(l)} = \alpha, \quad w_2^{(l)} = \frac{\alpha}{6}$$
 (9)

$$\mathcal{N}(z) = \frac{z}{\sqrt{z^{\dagger}z}} \tag{10}$$

$$\mathcal{G}_{\mu,n}(U^{(l)}) = V_{\mu}^{\dagger(l)}[U](n) \tag{11}$$

where $\mathcal{N}(z)$ is a projection or normalization function, which guarantees output of the layer is an element of the gauge group (or U(1) extended one with the normalization function). $V_{\mu}^{\dagger(l)}[U](n)$ is a staple.

2.2.2 Connection to stout smearing

Smearing step in the stout smearing (73) is expressed as

$$U_{\mu}^{(l+1)}(n) = \exp(Q_{\mu}^{(l)}(n))U_{\mu}^{(l)}(n)$$
(12)

$$= U_{\mu}^{(l)}(n) + \left(\exp(Q_{\mu}^{(l)}(n)) - 1\right) U_{\mu}^{(l)}(n). \tag{13}$$

So, we can identify,

$$w_1^{(l)} = w_2^{(l)} = 1 (14)$$

$$\mathcal{N}(z) = z \tag{15}$$

$$\mathcal{G}_{\mu,n}(U^{(l)}) = \left(\exp(Q_{\mu}^{(l)}(n)) - 1\right) U_{\mu}^{(l)}(n) \tag{16}$$

where $Q_{\mu}^{(l)}(n)$ is defined as

$$Q_{\mu}^{(l)}(n) = 2[\Omega_{\mu}^{(l)}(n)]_{\text{TA}}$$
(17)

and $\Omega_{\mu}^{(l)}(n)$ is constracted by untraced plaquette. Here TA means traceless-antihermitian operation.

Throughout this paper, we call the covariant neural network with Eq. (10) (11) as APE-type neural network, and with Eq. (14) (15) (16) (17) as stout-type neural network. As we will discuss later, the stout-type neural network can contain trainable parameters in $Q_{\mu}^{(l)}(n)$.

2.3 Rank-2 residual-net (Res-Net) and neural ordinary differential equation (Neural ODE)

2.3.1 Stout smearing and Rank-2 residual-net

We show that the covariant neural networks defined in Eq. (1) can be regarded as a rank-2 tensor version of the residual-net (Res-Net). If the covariant neural network is the stout type ($w_1 = 1$ and $\mathcal{N}(z) = z$), the tensors on the l-th layer is defined by (see e.g. Fig. 4)

$$U_{\mu}^{(l)}(n) = w_2^{(l-1)} \mathcal{G}_{\mu,n}^{\bar{\theta}^{(l-1)}}(U^{(l-1)}) + U_{\mu}^{(l-1)}(n)$$
(18)

This network structure is known as the Res-Net if U is a set of scalar values (Fig. 5). Thus, we call Eq. (18) the rank-2 Res-Net whose component is a rank-2 tensor. In the stout smearing scheme in the Lattice QCD community, there is no trainable parameter. However, if we regard the stout type smearing as the rank-2 Res-Net, we can introduce trainable parameters and find a back propagation of loss functions.

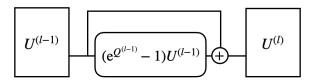


Figure 4: Calculation graph notation of neural networks for stout smearing. Input and output of this component are matrices. It is constructed by exponential function of un-traced plaquette.

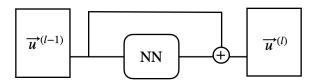


Figure 5: Calculation graph notation of Res-Net. Input and output of this component are vectors.

2.3.2 Rank-2 Neural ODE and gradient flow

The neural ODE is a cutting-edge framework of the deep learning. It is well known that the neural ODE is a generalization of the scalar-version of the Res-Net. If the layer index l is regarded as the discrete fictitious time t, the rank-2 neural ODE is defined as

$$\frac{dU_{\mu}^{(t)}(n)}{dt} = w_2^{(t)} \mathcal{G}_{\mu,n}^{\bar{\theta}^{(t)}}(U^{(t)}). \tag{19}$$

One can derive the rank-2 Res-Net (18) by solving the above differential equation with the use of the Euler method $(y_{n+1} = y_n + hf(t_n, y_n)$, for small h). On the other hand, the stout smearing is continuous version of the gradient flow with Wilson plaquette action and t is regarded as the flow time. If $\mathcal{G}_{\mu,n}^{\bar{\theta}^{(t)}}(U^{(t)})$ in equation above is taken as gradient of a gauge action, it is Luscher's gradient flow. Thus, we point out that the gradient flow with given parameters is regarded as the neural ODE for a gauge field in terms of the machine learning technique.

2.4 Rank-2 back propagation

Weight parameters in a neural network can be tuned by minimizing a loss function, which can be done with a gradient optimizer. Optimization process can be symbolically written,

$$\theta \leftarrow \theta - \eta \frac{\partial L_{\theta}(\mathcal{D})}{\partial \theta},\tag{20}$$

where η is a learning rate (a hyper parameter). In the calculation, we need a gradient of the loss function with respect to each weight $\frac{\partial L_{\theta}(\mathcal{D})}{\partial \theta}$. The derivatives are calculated by the delta rule and

backpropagation in the machine-learning community. For example, a weight in the last layer is,

$$\frac{\partial L_{\theta}(\mathcal{D})}{\partial w_i^{(L-1)}} = \frac{\partial L_{\theta}(\mathcal{D})}{\partial S_{\theta}} \frac{\partial S_{\theta}}{\partial w_i^{(L-1)}}$$
(21)

The first factor in the right hand side is calculable if concrete form of the loss function is given while second factor requires derivative with respect to $z_{\mu}^{(L)}(n)$, which is a rank-2 object. Thus, we cannot use conventional delta rule and we have to develop a rank-2 version.

2.4.1 Chain rules and star product

The scalar version of the back propagation scheme is derived from so-called the delta rule (See Appendix A.1). In this section, we derive the rank-2 version of the delta rule.

We introduce the matrix-valued "local" derivatives $\delta_{\mu}^{(l)}(n)$ on l-th layer, whose matrix element is defined as

$$\left[\delta_{\mu}^{(l)}(n)\right]^{i}{}_{j} \equiv \left[\frac{\partial S_{\theta}}{\partial z_{\mu}^{(l)}(n)}\right]^{i}{}_{j} \tag{22}$$

$$[\bar{\delta}_{\mu}^{(l)}(n)]^{i}_{j} \equiv \left[\frac{\partial S_{\theta}}{\partial z_{\mu}^{(l)\dagger}(n)}\right]^{i}_{j} \tag{23}$$

Here, we use definition of matrix derivative on a scalar function f as

$$\left[\frac{\partial f}{\partial A}\right]^{i}{}_{j} \equiv \frac{\partial f}{\partial A^{j}{}_{i}} \tag{24}$$

as in $(74)^5$. With the use of $\delta_{\mu}^{(l)}(n)$, the derivative with respect to the parameter can be calculated with a chain rule for a matrix

$$\frac{\partial S_{\theta}}{\partial \theta^{(l-1)}} = \sum_{\mu,n} \sum_{ij} \frac{\partial S_{\theta}}{\partial [z_{\mu}^{(l)}(n)]^{i}_{j}} \frac{\partial [z_{\mu}^{(l)}(n)]^{i}_{j}}{\partial \theta^{(l-1)}} = \sum_{\mu,n} \operatorname{Tr} \left[\frac{\partial S_{\theta}}{\partial z_{\mu}^{(l)}(n)} \frac{\partial z_{\mu}^{(l)}(n)}{\partial \theta^{(l-1)}} \right], \tag{25}$$

where $\theta^{(l-1)}$ is a parameter in (l-1)-th layer. If $z_{\mu}(n)$ is a function respect to $A_{\mu,m}$, a matrix-valued derivative $\partial f/\partial z_{\mu}(n)$ can be also calculated with a chain rule

$$\frac{\partial f}{\partial z_{\mu}(n)} = \sum_{\mu',m} \frac{\partial f}{\partial A_{\mu',n'}} \star \frac{\partial A_{\mu',n'}}{\partial z_{\mu}(n)} \tag{26}$$

This is not what is used in derivation of delta rule for conventional neural network but final results become same since HMC force is defined as a total variation. See (74). This definition gives $\frac{\partial \operatorname{tr} AB}{\partial A} = B$ instead of B^{\top} .

Here, we introduce the star product for the rank-2 and rank-4 tensors defined as

$$[A \star T]^i{}_j \equiv \sum_{kl} A^l{}_k T^k{}_j{}^i{}_l, \tag{27}$$

and the rank-4 tensor $\partial M/\partial A$ is defined as

$$\left[\frac{\partial M}{\partial A}\right]_{j}^{i}{}_{l}^{k} \equiv \frac{\partial}{\partial B^{j}{}_{k}} [M]^{i}{}_{l} \tag{28}$$

By introducing the star product for the rank-4 and rank-4 tensors defined as

$$(S \star T)^{i}{}_{j}{}^{k}{}_{l} = \sum_{nm} S^{i}{}_{n}{}^{m}{}_{l} T^{n}{}_{j}{}^{k}{}_{m}, \tag{29}$$

we can use a tensor version of chain rules:

$$\frac{\partial f}{\partial A} \equiv \sum_{IJ} \frac{\partial f}{\partial B_I} \star \frac{\partial B_I}{\partial C_J} \star \frac{\partial C_J}{\partial A},\tag{30}$$

where I, J will be identified as a coordinate and direction indices. Here we used associativity of the star product.

2.4.2 Rank-2 Delta rule and HMC force

HMC is the *de facto* standard algorithm in lattice QCD because it can deal with dynamical fermions (*33*). HMC is based on the Metropolis algorithm and the molecular dynamics to obtain configurations. The molecular dynamics uses gradient of the action with respect to the gauge field, which is called HMC force. For smeared action, one has to treat smeared links as a composite function of un-smeared (=thin) links, which is calculated by the chain rule. The HMC force with the covariant neural network is defined as

$$F_{\mu,n}[U] = \left[U_{\mu}(n) \frac{\partial S[U^{(L)}]}{\partial U_{\mu}(n)} \right]_{\text{TA}} = \left[U_{\mu}(n) \delta_{\mu}^{(l=0)}(n) \right]_{\text{TA}}, \tag{31}$$

where $U^{(L)}$ is the output from the covariant neural network and L is the number of covariant layers. The HMC force can be written in $\delta_{\mu}^{(l=0)}(n)$, which actually same object with $\Sigma(n)$ in (73).

With the star products and tensor version of chain rules, we can straightforwardly derive rank-2 version of the delta rule, which is a recurrence formula of $\delta_{\mu}^{(l)}(n)$. We note that there were similar derivations of HMC force with specific smearing schemes. To derive the recurrence formula of $\delta_{\mu}^{(l)}(n)$, we use the following chain rule,

$$\frac{\partial S}{\partial z_{\mu}^{(l)}(n)} = \sum_{\mu',n'} \left[\frac{\partial S}{\partial z_{\mu'}^{(l+1)}(n')} \star \frac{\partial z_{\mu'}^{(l+1)}(n')}{\partial z_{\mu}(n)} + \frac{\partial S}{\partial z_{\mu'}^{(l+1)\dagger}(n')} \star \frac{\partial z_{\mu'}^{(l+1)\dagger}(n')}{\partial z_{\mu}(n)} \right]. \tag{32}$$

This is an extended Wiltinger derivative to complex matrix (74) and constraint of special unitary is taken care of.

On the final layer $l=L,\,\delta_{\mu}^{(L)}(n)$ becomes

$$\delta_{\mu}^{(L)}(n) = \sum_{\alpha = L, IJ} \frac{\partial S_{\theta}}{\partial U_{\mu}^{(L)\alpha}(n)} \star \frac{\partial \mathcal{N}[z_{\mu}^{(L)}(n)]^{\alpha}}{\partial z_{\mu}^{(L)}(n)},\tag{33}$$

Here, we define $A^I(n) \equiv A$ and $A^{II} \equiv A^{\dagger}$.

On the *l*-th layer, $\delta_{\mu}^{(l)\alpha}(n)$ is written as

$$\delta_{\mu,\alpha}^{(l)}(n) = \sum_{\mu',m,\beta} \delta_{\mu',\beta}^{(l+1)}(m) \star \frac{\partial z_{\mu'}^{(l+1)\beta}(m)}{\partial z_{\mu}^{(l)\alpha}(n)}, \qquad (34)$$

$$= \sum_{\beta=I,II} \left\{ w_1^{(l)} \delta_{\mu,\beta}^{(l+1)}(n) \star \frac{\partial \mathcal{N}[z_{\mu}^{(l)\beta}(n)]}{\partial z_{\mu}^{(l)\alpha}(n)} + w_2^{(l)} \sum_{\mu',m} \delta_{\mu',\beta}^{(l+1)}(m) \star \sum_{\gamma=I,II} \frac{\partial \mathcal{G}_{\mu',m}(U^{(l)})^{\beta}}{\partial U_{\mu}^{(l)\gamma}(n)} \star \frac{\partial \mathcal{N}[z_{\mu}^{(l)}(n)]^{\gamma}}{\partial z_{\mu}^{(l)\alpha}(n)} \right\}, \qquad (35)$$

where $\alpha, \gamma = I, II$. This is the rank-2 delta rule for the covariant neural network. And by using this delta rule, we can optimize the weight in the network through a gradient optimizer.

2.5 Parametrization for stout-type covariant neural network

For stout-type covariant neural network, one has to take $w_1^{(l)}=w_2^{(l)}=1$ otherwise $\mathcal N$ has to take as normalization function as n-HYP. Weights can be included in $\Omega_\mu^{(l)}$. For example we can choose as,

$$\Omega_{\mu}^{(l)}(n) = \rho_{\text{plaq}}^{(l)} O_{\mu}^{\text{plaq}}(n) + \rho_{\text{rect}}^{(l)} O_{\mu}^{\text{rect}}(n) + \rho_{\text{poly}}^{(l)} O_{\mu}^{\text{poly}}(n), \tag{36}$$

where $O_{\mu}^{\text{plaq}}(n)$, $O_{\mu}^{\text{rect}}(n)$, $O_{\mu}^{\text{poly}}(n)$, are plaquette, rectangular loop, and Polyakov loop operators for μ direction originated from a point n, respectively. Every operators in $\Omega_{\mu}^{(l)}(n)$ is not traced as in stout and HEX smearing.

We can derive a derivative with respect to ρ using formula above,

$$\frac{\partial S}{\partial \rho_i^{(l)}} = 2 \operatorname{Re} \sum_{\mu',m} \operatorname{tr} \left[U_{\mu'}^{(l)\dagger}(m) \Lambda_{\mu',m} \frac{\partial C}{\partial \rho_i^{(l)}} \right], \tag{37}$$

where $\rho_i^{(l)}$ is a coefficient of a loop operator type i (plaquette, rectangler, Polyakov loop, etc) in the level l. C is sum of staples for loop operators in $\Omega_{\mu}^{(l)}(n)$. $\Lambda_{\mu,m}$ is defined in (73).

2.6 Connection to fully connected neural network

The covariant neural network fallback to a conventional neural network a linear function of a gauge field $A_{\mu}(n)$ if a is small for $U_{\mu}(n) = \mathrm{e}^{\mathrm{i} a A_{\mu}(n)} \approx 1 + \mathrm{i} a A_{\mu}(n)$. We follow the idea in (75), perturbative analysis of smearing. For example, an argument of Eq. (2) Eq. (3) for APE type smearing with a staple for plaquette is,

$$w_{1}^{(l-1)}U_{\mu}(n) + w_{2}^{(l-1)} \sum_{\mu \neq \nu} [U_{\nu}(n+\hat{\mu})U_{\mu}^{\dagger}(n+\hat{\nu})U_{\nu}^{\dagger}(n) + \text{opposite}], \tag{38}$$

$$= w_{1}^{(l-1)}(1+iaA_{\mu}(n)) + w_{2}^{(l-1)} \sum_{\mu \neq \nu} [(1+iaA_{\nu}(n+\hat{\mu}))(1-iaA_{\mu}(n+\hat{\nu}))(1-iaA_{\nu}(n)) + (\text{opposite})], \tag{39}$$

$$= w_{1}^{(l-1)}iaA_{\mu}(n) + w_{2}^{(l-1)} \sum_{\mu \neq \nu} [iaA_{\nu}(n+\hat{\mu}) - iaA_{\mu}(n+\hat{\nu}) - iaA_{\nu}(n) + (\text{opposite})] + \text{const} + O(a^{2}), \tag{40}$$

where (opposite) indicates staple along with the opposite side. This is nothing but a linear summation of input. This is fed to an non-linear function, which corresponds to an activation function. Moreover, if we do not use weight sharing, the right hand side of equation above is

extended as,

$$w_{1,\mu}^{(l-1)}(n)iaA_{\mu}(n) + \sum_{\mu \neq \nu} [iaw_{2,\nu}^{(l-1)}(n+\hat{\mu})A_{\nu}(n+\hat{\mu}) - iaw_{2,\nu}^{(l-1)}(n+\hat{\nu})A_{\mu}(n+\hat{\nu}) - iaw_{2,\nu}^{(l-1)}(n)A_{\nu}(n) + (\text{opposite})] + \text{const} + O(a^{2}).$$
(41)

If the gauge group is U(1), the gauge field $A_{\mu}(n)$ becomes real valued. This can be regarded as a dense layer in a conventional neural network if we identify a pair of index (μ, n) as index for a input vector in conventional fully connected layer.

3 Demonstration

Here we show how the gauge covariant neural network works using self-learning hybrid Monte-Carlo (SLHMC) (64) as a demonstration. SLHMC uses the target action in the Metropolis test and machine learning approximated action in the molecular dynamics in HMC and we use gauge covariant neural network approximated action.

3.1 Self-learning Hybrid Monte-Carlo for QCD

SLHMC is an exact algorithm, which is consisted by two parts, the molecular dynamics and the Metropolis test. We use S[U], action for the target system, in the Metropolis test and $S_{\theta}^{\text{eff}}[U]$, an effective action for the system with a set of parameters θ , in the molecular dynamics. Since the molecular dynamics is reversible, we do not have to perform the Metropolis-Hasting test. Instead, the Metropolis test is enough to guarantee the convergence (53,64). The same mechanism has been used in the rational hybrid Monte-Carlo (34).

Our target system is written as,

$$S[U] = S_{g}[U] + S_{f}[\phi, U; m_{l}], \tag{42}$$

where $S_{\mathbf{g}}[U]$ is a gauge action and $S_{\mathbf{f}}[\phi,U;m_{\mathbf{l}}]$ is an action for a psuedo-fermion field ϕ with

mass m_1 . In HMC, total action is a function of a pseudo-fermion field but here we suppress in equation for notational simplicity.

In SLHMC for QCD, we can take following effective action in general,

$$S_{\theta}[U] = S_{\theta}^{\text{eff}}[U_{\theta'}^{\text{NN}}[U]] + S_{\text{f}}^{\text{eff}}[\phi, U_{\theta}^{\text{NN}}[U]; m_{\text{h}}], \tag{43}$$

where $U_{\theta}^{\rm NN}[U]$ is a transformed configuration with a trained covariant neural network and $U_{\theta'}^{\rm NN}[U]$ can be another network. $m_{\rm h}$ is fermion mass for effective action, which will be taken as $m_{\rm h} > m_{\rm l}$. $S_{\rm g}^{\rm eff}$ is a gauge action, which can contain tunable parameter as in (53) but here we just use Wilson plaquette action. $S_{\rm f}^{\rm eff}$ is a fermion action again and it can be different with the target one and we will mention this extension but we use the same action with the target. Namely we take,

$$S_{\theta}[U] = S_{g}[U] + S_{f}[\phi, U_{\theta}^{NN}[U]; m_{h}], \tag{44}$$

in our calculation. This is a neural network parametrized effective action. If the neural network has enough expressibility, as (3-7), we expect that the effective action can enough close to the target one as an functional of U and ϕ .

We use stout type covariant neural network Eq. (16) and Eq. (17) with

$$\Omega_{\mu}^{(l)}(n) = \rho_{\text{plaq}}^{(l)} O_{\mu}^{\text{plaq}}(n) + \begin{cases} \rho_{\text{poly,4}}^{(l)} O_{4}^{\text{poly}}(n) & (\mu = 4), \\ \rho_{\text{poly,s}}^{(l)} O_{i}^{\text{poly}}(n), & (\mu = i = 1, 2, 3) \end{cases}$$
(45)

and we omit rectangular term to reduce numerical cost and for simplicity. In our previous work (53), we had to use fermion determinant directory since SLMC for QCD requires absolute value of fermionic determinant in the training and production. However in the current work, we do not have to use the determinant because it is relay on SLHMC framework. Note that, our choice Eq. (45) is an natural extension to an effective action our previous work (53).

3.1.1 Training

Here we describe how to train our neural network. We train S_{θ} in HMC run prior to SLHMC run. After the Metropolis step in HMC, we keep the pseudo-fermion field and a gauge configuration U, and we calculate,

$$L_{\theta}[U] = \frac{1}{2} |S_{\theta}[U, \phi] - S[U, \phi]|^2,$$
 (46)

and we perform the backpropagation with the delta rule. This is a gauge invariant loss function. We produce gauge configurations with trained neural network after the training for simplicity. Training and production can be done simultaneously as in the previous paper (53).

3.1.2 Configuration generation with SLHMC

After the training, we generate configurations using the neural network. As we stated, we use the target action S in the Metropolis test in SLHMC, and S_{θ} in the molecular dynamics. We fix the weights during the production run for simplicity.

3.2 Numerical setup

Here we summarize our numerical setup. We implement⁶ our framework for SU(2) gauge theory in four dimension with dynamical fermions. Our code is implemented based on Lattice-QCD.jl (76) with Julia (77). We use automatic staple generator from loops operators in the delta rule, which has been used in our public code (76). In addition we implement automatic staple derivative generator for given general staples.

We perform simulation with staggered fermions without rooting and plaquette gauge action in $N_{\sigma}^3 N_{\tau} = 4^4$ with $m_1 = 0.3$ and $\beta = 2.7$ for this proof-of-principle study. We train our network with $m_{\rm h} = 0.4, 0.5, 0.75, 1.0$ (Tab. 1) to see a training history of the loss function and weights.

⁶Our framework can be implemented immediately if a public code has smearing functionality.

We choose our training rate $\eta = 1.0 \times 10^{-7}$, which is not well tuned but enough for this demonstration and employ the simplest stochastic gradient as an optimizer. Training was performed for 1500 trajectories (training step) from a thermalized configuration.

The SLHMC run performed 4500 trajectories with trained weights and $m_h = 0.4$, and first 200 trajectories are not included in the analysis. The step size in the molecular dynamics is taken to 0.02 to see quality of trained action.

 N_{τ}	N_{σ}	m_{h}	$N_{ m layers}$	Loops
4	4	1.0	2	Plaquette, Polyakov loops (spatial, temporal)
4	4	0.75	2	Plaquette, Polyakov loops (spatial, temporal)
4	4	0.5	2	Plaquette, Polyakov loops (spatial, temporal)
4	4	0.4	2	Plaquette, Polyakov loops (spatial, temporal)

Table 1: "Loops" means types of untraced loop operator in the stout-type covariant neural networks. In each case, the network contains 6 parameters.

3.3 Results

Here we show our results. Left panel of Fig. 6 is the history of the loss function for various $m_{\rm h}$ in the log scale along with the molecular dynamics time in the prior HMC, which coincides with the training steps in the current case. $m_{\rm h}=1.0$ is not converging around the molecular dynamics time around 1500, we train it for more than 10000 steps. For $m_{\rm l}\sim m_{\rm h}$ case, the loss functions fluctuate around 1. For $m_{\rm h}\gg m_{\rm l}$, it starts from $\sim O(10^3)$ and it decreases to O(1). These results indicate that, the neural network is correctly trained with the stochastic gradient and our formulae. Right panel of Fig. 6 is loss history for $m_{\rm h}=0.4$ in linear the scale. In side plot is zoomed around MD time equals to 0. For earlier stage of training, it starts from 89, and it decreases with fluctuation. After some steps, it fluctuates O(1).

Fig. 7 shows training history of weights in the neural network in HMC. From top to bottom, coefficients of plaquette, Polyakov loop for spatial directions, and Polyakov loop for temporal directions are displayed. Left panels and right panels are weights for the first layer and the

second layer. One can see that all of weights $\rho \neq 0$ for after the training. All of weights in the first layer show same tendency with ones in the second layer. This means, some of weights are irrelevant and we can reduce the degrees of freedom by adding a regulator term into the loss function. Coefficients of Polyakov loop for spatial direction has opposite tendency with temporal ones. the weight for Polyakov loop with $m_h = 1.0$ at l = 1 has short plateau MD time = 1200 and it goes down again around 1800. This indicates that, if we employ an effective action far from the target action, longer training steps are needed. In SLMC run, we use trained weight in 1500-th step (Tab. 2). Some of values are negative, while coefficients in conventional smearing are positive.

Layer	Loop	Value of ρ
1	Plaquette	-0.011146476388409423
2	Plaquette	-0.011164492428633698
1	Spatial Polyakov loop	-0.0030283193221172216
2	Spatial Polyakov loop	-0.0029984533773388094
1	Temporal Polyakov loop	0.004248021727233112
2	Temporal Polyakov loop	0.004195253380373369

Table 2: Trained weights for the current setup. This value is used in the SLHMC run.

Fig. 8 shows histogram of observables of HMC and SLHMC with $m_{\rm h}=0.4$. Top left is plaquette, top right is Polyakov loop for temporal direction, and bottom is the chiral condensate, respectively. Error bars are evaluated by the Jackknife method. All of histograms are well overlapped as we expected. Tab. 3 is summary table for results from HMC and SLMC. Autocorrelation effects in statistical error are taken into account in statistical error evaluation. Two algorithms give consistent results in error.

Finally, we comment on the acceptance rate of this demonstration. Acceptance rate of SLHMC with $m_{\rm h}=0.4$ in this 6 parameters run is roughly 40%. This seems not efficient but it can be improved by employing more expressible neural network adding more loops and layers. Owing to the training, we achieve non-zero acceptance with a simple parametrized

action with different mass.

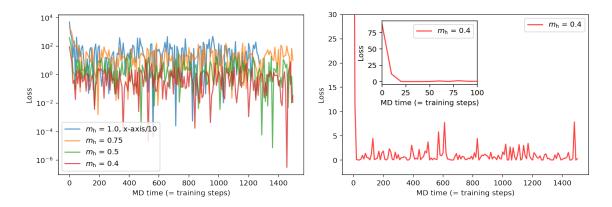


Figure 6: Loss history in prior run of HMC. Horizontal axis coincides with the training steps. (*Left*) Summary plot for $m_h = 1.0, 0.75, 0.5, 0.4$ in the log scale. (*Right*) Loss history for $m_h = 0.4$ in the linear scale. The inside plot is zoomed around MD time equals to zero.

Algorithm	Observable	Value
HMC	Plaquette	0.7024(2)
SLHMC	Plaquette	0.7021(4)
HMC	Polyakov loop	0.84(2)
SLHMC	Polyakov loop	0.82(3)
HMC	Chiral condensate	0.4238(5)
SLHMC	Chiral condensate	0.4261(10)

Table 3: Summary table for results from HMC and SLMC.

4 Summary and discussion

In this paper, we introduced a neural network, which make a map between gauge configurations with gauge covariance. This neural network can be regarded as smearing with tunable parameters. This is analogous to the convolutional layer, which is regarded as tunable image filter. Linear transformation in conventional neural networks corresponds to summation of links and staples in smearing and non-linear transformation for APE-type case. For stout-type, it is summation of un-traced loop operators in the exponential function. Non-linear transformation in neural networks corresponds to projection or normalization process in APE type smearing. For stout type, exponential function plays a role for the activation function. Using the output of the network, we can define real valued gauge invariant loss function through gauge invariant actions.

As a use case of the covariant neural network, we performed simulation of SLHMC (64) for two color QCD with un-rooted staggered fermion, which corresponds to 4 dynamical fermions in the continuum. In contrast to the previous work (53), we can perform simulation without calculating the fermion determinant. We parametrized fermionic action with the neural network and observed decreasing of the loss function in the training. In the production phase, we obtained consistent results to HMC.

Here we enumerate possible applications. One can perform systematic search of neural network effective action which is able to achieve higher acceptance than the demonstration in this paper. Main focus of this paper is to develop the covariant neural network, so this would be our future work.

If we take S_{θ} as the standard staggered fermion action and fix the parameters in the neural network, it becomes a smeared fermions like stout smeared staggered fermions. If we choose APE-like smearing with extended links (up to 7 links) with the staggered Dirac operator and appropriately tuned weight and parameters, and then, we essentially get ASQTAD (A-squared and tad-pole improved staggered) quarks and HISQ (highly improved staggered quarks) (78–80) can be obtained similar manner⁷. In our framework, we could construct much complicated and better action using training.

As we have shown that, neural network parametrized SLHMC works with lattice fermions in contrast to our previous work (53). We may be able to construct overlap like domain-wall

⁷HISQ contains both of level-3 and level-2 smeared links, thus construction of network becomes complected but still HMC force can be obtained using the delta rule.

fermions with a covariant neural network. Namely we take $S_{\rm f}$ as overlap fermion action and $S_{\rm f}^{\rm eff}$ as the action for domain-wall fermion with $U_{\theta}^{\rm NN}[U]$.

Using covariant neural networks, we can realizes gauge invariant ultraviolet filtering to remove ultraviolet noise in the topological charge as in the conventional gradient flow. By using established techniques in neural ODE (24), we can design better ultraviolet filter with shorter flow time. HMC with an action $Q_{\text{top}}[U_{\theta}^{\text{NN}}[U]]$ as (81).

If one takes weight w_k^l as complex, this can realize an anti-holomorphic flow like (82). By using covariant neural net, one can extend their framework in (82) to non-abelian gauge theory in higher dimension.

Topology conservation in neural net has been discussed (83, 84). One could prove conservation of topology in the covariant neural networks.

Beyond the lattice QCD and physics, it might be applicable to general machine learning. Taco Cohen and Max Welling *et al.* invented gauge equivalent neural networks to for data with descrete gauge symmetry (17). Our framework guarantees gauge covariance for continuous group, thus, some data with local continuous symmetry. If one implement recognition architecture with covariant neural network, performance would be improved.

In our framework for SLHMC for QCD, we can exact loss function because the probability distribution is known unlike the other machine learning dataset. One could calculate the loss landscape using perturbation or strong coupling expansion for lattice QCD.

The current architecture uses non-linear functions as in the conventional neural networks, which is potentially a universal approximator and the property is supported by the non-linear functions and stacking of many layers as in other deep learning architectures (4–7). It is important to proof, the covariant neural network is a universal approximator for gauge covariant local functions but this is out of scope of this work and we leave this for future work.

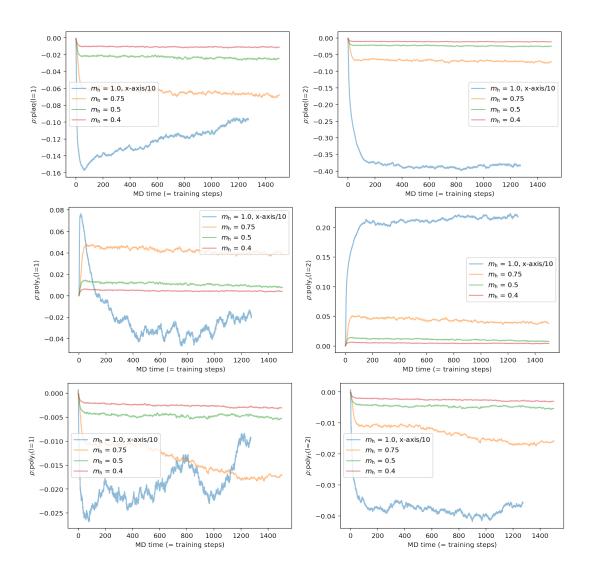


Figure 7: Training history of weights in the neural network. From top to bottom, coefficients of plaquette, Polyakov loop for spatial directions, and Polyakov loop for temporal directions are displayed. Left panels and right panels are weights for the first layer and the second layer.

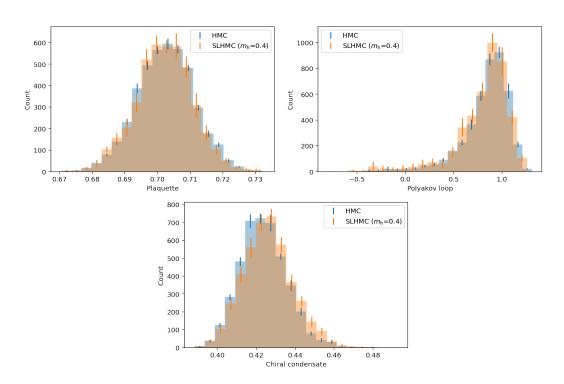


Figure 8: Comparison of results for HMC and SLHMC with $m_{\rm h}=0.4$. Top left is plaquette, top right is Polyakov loop for temporal direction, and bottom is the chiral condensate, respectively.

References

- 1. Giuseppe Carleo, Ignacio Cirac, Kyle Cranmer, Laurent Daudet, Maria Schuld, Naftali Tishby, Leslie Vogt-Maranto, and Lenka Zdeborová. Machine learning and the physical sciences. *Reviews of Modern Physics*, 91(4), Dec 2019.
- Pankaj Mehta, Marin Bukov, Ching-Hao Wang, Alexandre G.R. Day, Clint Richardson, Charles K. Fisher, and David J. Schwab. A high-bias, low-variance introduction to machine learning for physicists. *Physics Reports*, 810:1–124, May 2019.
- 3. George Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of control, signals and systems*, 2(4):303–314, 1989.
- 4. Ding-Xuan Zhou. Universality of deep convolutional neural networks, 2018.
- 5. Patrick Kidger and Terry Lyons. Universal approximation with deep narrow networks, 2020.
- 6. Paulo Tabuada and Bahman Gharesifard. Universal approximation power of deep residual neural networks via nonlinear control theory, 2020.
- 7. Jesse Johnson. Deep, skinny neural networks are not universal approximators, 2018.
- 8. Lenka Zdeborová. New tool in the box. *Nature Physics*, 13(5):420–421, 2017.
- 9. Juan Carrasquilla and Roger G Melko. Machine learning phases of matter. *Nature Physics*, 13(5):431–434, 2017.
- 10. Giuseppe Carleo and Matthias Troyer. Solving the quantum many-body problem with artificial neural networks. *Science*, 355(6325):602–606, 2017.

- 11. Akinori Tanaka and Akio Tomiya. Detection of phase transition via convolutional neural networks. *Journal of the Physical Society of Japan*, 86(6):063001, 2017.
- 12. Kouji Kashiwa, Yuta Kikuchi, and Akio Tomiya. Phase transition encoded in neural network. *Progress of Theoretical and Experimental Physics*, 2019(8):083A04, 2019.
- 13. M. Z. Hasan and C. L. Kane. Colloquium: Topological insulators. *Reviews of Modern Physics*, 82(4):3045–3067, Nov 2010.
- 14. Particle Data Group, P A Zyla, R M Barnett, J Beringer, O Dahl, D A Dwyer, D E Groom, C J Lin, K S Lugovsky, E Pianori, D J Robinson, C G Wohl, W M Yao, K Agashe, G Aielli, B C Allanach, C Amsler, M Antonelli, E C Aschenauer, D M Asner, H Baer, Sw Banerjee, L Baudis, C W Bauer, J J Beatty, V I Belousov, S Bethke, A Bettini, O Biebel, K M Black, E Blucher, O Buchmuller, V Burkert, M A Bychkov, R N Cahn, M Carena, A Ceccucci, A Cerri, D Chakraborty, R Sekhar Chivukula, G Cowan, G D'Ambrosio, T Damour, D de Florian, A de Gouvêa, T DeGrand, P de Jong, G Dissertori, B A Dobrescu, M D'Onofrio, M Doser, M Drees, H K Dreiner, P Eerola, U Egede, S Eidelman, J Ellis, J Erler, V V Ezhela, W Fetscher, B D Fields, B Foster, A Freitas, H Gallagher, L Garren, H J Gerber, G Gerbier, T Gershon, Y Gershtein, T Gherghetta, A A Godizov, M C Gonzalez-Garcia, M Goodman, C Grab, A V Gritsan, C Grojean, M Grünewald, A Gurtu, T Gutsche, H E Haber, C Hanhart, S Hashimoto, Y Hayato, A Hebecker, S Heinemeyer, B Heltsley, J J Hernández-Rey, K Hikasa, J Hisano, A Höcker, J Holder, A Holtkamp, J Huston, T Hyodo, K F Johnson, M Kado, M Karliner, U F Katz, M Kenzie, V A Khoze, S R Klein, E Klempt, R V Kowalewski, F Krauss, M Kreps, B Krusche, Y Kwon, O Lahav, J Laiho, L P Lellouch, J Lesgourgues, A R Liddle, Z Ligeti, C Lippmann, T M Liss, L Littenberg, C Lourengo, S B Lugovsky, A Lusiani, Y Makida, F Maltoni, T Mannel, A V Manohar, W J Marciano, A Masoni, J Matthews, U G Meißner, M Mikhasenko, D J Miller, D Milstead, R E Mitchell,

K Mönig, P Molaro, F Moortgat, M Moskovic, K Nakamura, M Narain, P Nason, S Navas, M Neubert, P Nevski, Y Nir, K A Olive, C Patrignani, J A Peacock, S T Petcov, V A Petrov, A Pich, A Piepke, A Pomarol, S Profumo, A Quadt, K Rabbertz, J Rademacker, G Raffelt, H Ramani, M Ramsey-Musolf, B N Ratcliff, P Richardson, A Ringwald, S Roesler, S Rolli, A Romaniouk, L J Rosenberg, J L Rosner, G Rybka, M Ryskin, R A Ryutin, Y Sakai, G P Salam, S Sarkar, F Sauli, O Schneider, K Scholberg, A J Schwartz, J Schwiening, D Scott, V Sharma, S R Sharpe, T Shutt, M Silari, T Sjöstrand, P Skands, T Skwarnicki, G F Smoot, A Soffer, M S Sozzi, S Spanier, C Spiering, A Stahl, S L Stone, Y Sumino, T Sumiyoshi, M J Syphers, F Takahashi, M Tanabashi, J Tanaka, M Taševský, K Terashi, J Terning, U Thoma, R S Thorne, L Tiator, M Titov, N P Tkachenko, D R Tovey, K Trabelsi, P Urquijo, G Valencia, R Van de Water, N Varelas, G Venanzoni, L Verde, M G Vincter, P Vogel, W Vogelsang, A Vogt, V Vorobyev, S P Wakely, W Walkowiak, C W Walter, D Wands, M O Wascko, D H Weinberg, E J Weinberg, M White, L R Wiencke, S Willocq, C L Woody, R L Workman, M Yokoyama, R Yoshida, G Zanderighi, G P Zeller, O V Zenin, R Y Zhu, S L Zhu, F Zimmermann, J Anderson, T Basaglia, V S Lugovsky, P Schaffner, and W Zheng. Review of Particle Physics. Progress of Theoretical and Experimental Physics, 2020(8), 08 2020. 083C01.

- 15. Taco S. Cohen and Max Welling. Group equivariant convolutional networks, 2016.
- 16. Taco S. Cohen, Mario Geiger, Jonas Koehler, and Max Welling. Spherical cnns, 2018.
- 17. Taco S. Cohen, Maurice Weiler, Berkay Kicanaoglu, and Max Welling. Gauge equivariant convolutional networks and the icosahedral cnn, 2019.
- 18. Martin Lüscher. Trivializing maps, the wilson flow and the hmc algorithm. *Communications in Mathematical Physics*, 293(3):899–919, Nov 2009.

- 19. Cédric Bény. Deep learning and the renormalization group, 2013.
- 20. Henry W. Lin, Max Tegmark, and David Rolnick. Why does deep and cheap learning work so well? *Journal of Statistical Physics*, 168(6):1223–1247, Jul 2017.
- 21. Koji Hashimoto, Sotaro Sugishita, Akinori Tanaka, and Akio Tomiya. Deep learning and the ads/cft correspondence. *Physical Review D*, 98(4), Aug 2018.
- 22. Koji Hashimoto, Sotaro Sugishita, Akinori Tanaka, and Akio Tomiya. Deep learning and holographic qcd. *Physical Review D*, 98(10), Nov 2018.
- 23. Koji Hashimoto, Hong-Ye Hu, and Yi-Zhuang You. Neural ode and holographic qcd, 2020.
- 24. Ricky T. Q. Chen, Yulia Rubanova, Jesse Bettencourt, and David Duvenaud. Neural ordinary differential equations, 2019.
- Hong-Ye Hu, Shuo-Hui Li, Lei Wang, and Yi-Zhuang You. Machine learning holographic mapping by neural network renormalization group. *Physical Review Research*, 2(2), Jun 2020.
- 26. Sinya Aoki and Shuichi Yokoyama. Flow equation, conformal symmetry, and anti-de sitter geometry. *Progress of Theoretical and Experimental Physics*, 2018(3), Mar 2018.
- 27. Sinya Aoki and Shuichi Yokoyama. Ads geometry from cft on a general conformally flat manifold. *Nuclear Physics B*, 933:262–274, Aug 2018.
- 28. Vasudev Shyam. Finite cutoff ads5 holography and the generalized gradient flow. *Journal of High Energy Physics*, 2018(12), Dec 2018.
- 29. Sinya Aoki, Shuichi Yokoyama, and Kentaroh Yoshida. Holographic geometry for nonrelativistic systems emerging from generalized flow equations. *Physical Review D*, 99(12), Jun 2019.

- 30. Kenneth G. Wilson. Confinement of quarks. Phys. Rev. D, 10:2445–2459, Oct 1974.
- 31. Michael Creutz. Monte carlo study of quantized su(2) gauge theory. *Phys. Rev. D*, 21:2308–2315, Apr 1980.
- 32. Peter A Boyle. Machines and algorithms, 2017.
- 33. S. Duane, A. D. Kennedy, B. J. Pendleton, and D. Roweth. Hybrid Monte Carlo. *Phys. Lett. B*, 195:216–222, 1987.
- 34. M. A. Clark. The Rational Hybrid Monte Carlo Algorithm. *PoS*, LAT2006:004, 2006.
- 35. Kohtaroh Miura. Review of Lattice QCD Studies of Hadronic Vacuum Polarization Contribution to Muon g-2. *PoS*, LATTICE2018:010, 2019.
- 36. S. Aoki et al. FLAG Review 2019: Flavour Lattice Averaging Group (FLAG). *Eur. Phys. J. C*, 80(2):113, 2020.
- 37. Heng-Tong Ding, Frithjof Karsch, and Swagato Mukherjee. Thermodynamics of strong-interaction matter from Lattice QCD. *Int. J. Mod. Phys. E*, 24(10):1530007, 2015.
- 38. Martin Lüscher and Peter Weisz. Perturbative analysis of the gradient flow in non-abelian gauge theories. *Journal of High Energy Physics*, 2011(2), Feb 2011.
- 39. Martin Lüscher. Future applications of the yang-mills gradient flow in lattice qcd, 2013.
- 40. Aya Kagimura, Akio Tomiya, and Ryo Yamamura. Effective lattice action for the configurations smeared by the wilson flow, 2015.
- 41. Andrea Carosso, Anna Hasenfratz, and Ethan T. Neil. Stochastic renormalization group and gradient flow in scalar field theory, 2019.

- 42. Hidenori Sonoda and Hiroshi Suzuki. Derivation of a gradient flow from the exact renormalization group. *PTEP*, 2019(3):033B05, 2019.
- 43. Hidenori Sonoda and Hiroshi Suzuki. Gradient flow exact renormalization group. *Progress of Theoretical and Experimental Physics*, 2021(2), Jan 2021.
- 44. Sebastian J. Wetzel and Manuel Scherzer. Machine learning of explicit order parameters: From the ising model to su(2) lattice gauge theory. *Physical Review B*, 96(18), Nov 2017.
- 45. Stefan Blücher, Lukas Kades, Jan M. Pawlowski, Nils Strodthoff, and Julian M. Urban. Towards novel insights in lattice field theory with explainable machine learning. *Physical Review D*, 101(9), May 2020.
- 46. D. L. Boyda, M. N. Chernodub, N. V. Gerasimeniuk, V. A. Goy, S D. Liubimov, and A. V. Molochkov. Finding the deconfinement temperature in lattice yang-mills theories from outside the scaling window with machine learning. *Physical Review D*, 103(1), Jan 2021.
- 47. Boram Yoon, Tanmoy Bhattacharya, and Rajan Gupta. Machine learning estimators for lattice qcd observables. *Physical Review D*, 100(1), Jul 2019.
- 48. Rui Zhang, Zhouyou Fan, Ruizi Li, Huey-Wen Lin, and Boram Yoon. Machine-learning prediction for quasiparton distribution function matrix elements. *Physical Review D*, 101(3), Feb 2020.
- 49. Takuya Matsumoto, Masakiyo Kitazawa, and Yasuhiro Kohno. Classifying topological charge in SU(3) Yang–Mills theory with machine learning. *PTEP*, 2021(2):023D01, 2021.
- 50. Akinori Tanaka and Akio Tomiya. Towards reduction of autocorrelation in HMC by machine learning. 12 2017.

- 51. Kai Zhou, Gergely Endrődi, Long-Gang Pang, and Horst Stöcker. Regressive and generative neural networks for scalar field theory. *Phys. Rev. D*, 100(1):011501, 2019.
- 52. Jan M. Pawlowski and Julian M. Urban. Reducing Autocorrelation Times in Lattice Simulations with Generative Adversarial Networks. *Mach. Learn. Sci. Tech.*, 1:045011, 2020.
- 53. Yuki Nagai, Akinori Tanaka, and Akio Tomiya. Self-learning Monte-Carlo for non-abelian gauge theory with dynamical fermions. 10 2020.
- 54. M. S. Albergo, G. Kanwar, and P. E. Shanahan. Flow-based generative models for Markov chain Monte Carlo in lattice field theory. *Phys. Rev. D*, 100(3):034515, 2019.
- 55. Danilo Jimenez Rezende, George Papamakarios, Sébastien Racanière, Michael S. Albergo, Gurtej Kanwar, Phiala E. Shanahan, and Kyle Cranmer. Normalizing flows on tori and spheres, 2020.
- 56. Gurtej Kanwar, Michael S. Albergo, Denis Boyda, Kyle Cranmer, Daniel C. Hackett, Sébastien Racanière, Danilo Jimenez Rezende, and Phiala E. Shanahan. Equivariant flow-based sampling for lattice gauge theory. *Phys. Rev. Lett.*, 125(12):121601, 2020.
- 57. Denis Boyda, Gurtej Kanwar, Sébastien Racanière, Danilo Jimenez Rezende, Michael S. Albergo, Kyle Cranmer, Daniel C. Hackett, and Phiala E. Shanahan. Sampling using SU(N) gauge equivariant flows. 8 2020.
- 58. Michael S. Albergo, Denis Boyda, Daniel C. Hackett, Gurtej Kanwar, Kyle Cranmer, Sébastien Racanière, Danilo Jimenez Rezende, and Phiala E. Shanahan. Introduction to Normalizing Flows for Lattice Field Theory. 1 2021.
- 59. Shuo-Hui Li and Lei Wang. Neural network renormalization group. *Physical Review Letters*, 121(26), Dec 2018.

- 60. James Halverson, Anindita Maiti, and Keegan Stoner. Neural networks and quantum field theory, 2020.
- 61. Koji Hashimoto. AdS/CFT correspondence as a deep Boltzmann machine. *Phys. Rev. D*, 99(10):106017, 2019.
- 62. Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015.
- 63. Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks, 2016.
- 64. Yuki Nagai, Masahiko Okumura, Keita Kobayashi, and Motoyuki Shiga. Self-learning hybrid monte carlo: A first-principles approach. *Physical Review B*, 102(4), Jul 2020.
- 65. Di Luo, Giuseppe Carleo, Bryan K. Clark, and James Stokes. Gauge equivariant neural networks for quantum lattice gauge theories. 12 2020.
- 66. Danilo Jimenez Rezende, Sébastien Racanière, Irina Higgins, and Peter Toth. Equivariant hamiltonian flows, 2019.
- 67. S. Groote, R. Saar, and H. Liivat. Lorentz invariance and gauge equivariance. *J. Phys. Conf. Ser.*, 532:012007, 2014.
- 68. Diego Marcos, Michele Volpi, Nikos Komodakis, and Devis Tuia. Rotation equivariant vector field networks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5048–5057, 2017.
- 69. Jörg Behler and Michele Parrinello. Generalized neural-network representation of high-dimensional potential-energy surfaces. *Phys. Rev. Lett.*, 98:146401, Apr 2007.

- 70. Yuki Nagai, Masahiko Okumura, and Akinori Tanaka. Self-learning monte carlo method with behler-parrinello neural networks. *Physical Review B*, 101(11), Mar 2020.
- 71. M. Albanese et al. Glueball Masses and String Tension in Lattice QCD. *Phys. Lett. B*, 192:163–169, 1987.
- 72. Anna Hasenfratz, Roland Hoffmann, and Stefan Schaefer. Hypercubic smeared links for dynamical fermions. *Journal of High Energy Physics*, 2007(05):029–029, May 2007.
- 73. Colin Morningstar and Mike Peardon. Analytic smearing ofsu(3)link variables in lattice qcd. *Physical Review D*, 69(5), Mar 2004.
- 74. Waseem Kamleh, Derek B. Leinweber, and Anthony G. Williams. Hybrid Monte Carlo with fat link fermion actions. *Phys. Rev. D*, 70:014502, 2004.
- 75. Stefano Capitani, Stephan Durr, and Christian Hoelbling. Rationale for UV-filtered clover fermions. *JHEP*, 11:028, 2006.
- 76. Yuki Nagai and Akio Tomiya. Latticeqcd.jl. Paper is in preparation, https://github.com/akio-tomiya/LatticeQCD.jl.
- 77. Jeff Bezanson, Alan Edelman, Stefan Karpinski, and Viral B. Shah. Julia: A fresh approach to numerical computing, 2015.
- 78. G. Peter Lepage. Flavor symmetry restoration and Symanzik improvement for staggered quarks. *Phys. Rev. D*, 59:074502, 1999.
- 79. Kostas Orginos, Doug Toussaint, and R. L. Sugar. Variants of fattening and flavor symmetry restoration. *Phys. Rev. D*, 60:054503, 1999.

- 80. Kostas Orginos and Doug Toussaint. Testing improved actions for dynamical Kogut-Susskind quarks. *Phys. Rev. D*, 59:014501, 1999.
- 81. Claudio Bonati, Massimo D'Elia, Guido Martinelli, Francesco Negro, Francesco Sanfilippo, and Antonino Todaro. Topology in full QCD at high temperature: a multicanonical approach. *JHEP*, 11:170, 2018.
- 82. Andrei Alexandru, Paulo F. Bedaque, Henry Lamm, and Scott Lawrence. Deep learning beyond lefschetz thimbles. *Physical Review D*, 96(9), Nov 2017.
- 83. Mustafa Hajij and Kyle Istvan. A topological framework for deep learning, 2021.
- 84. C. Olah. Neural Networks, Manifolds, and Topology. 2021.
- 85. Koji Hashimoto Akinori Tanaka, Akio Tomiya. *Deep Learning and Physics*. Springer, Switzerland, 2021.

5 Acknowledgments

The work of A.T. was supported by the RIKEN Special Postdoctoral Researcher program. YN was partially supported by JSPS- KAKENHI Grant Numbers 18K11345 and 18K03552.

A Appendix

A.1 Conventional Neural networks

Here we introduce conventional Neural networks to be self-contained. Please refer (85) to see detail for example. Let us take $\vec{x} \in \mathbb{R}^n$ is an input vector of a neural network. A neural network is symbolically written as, a composite function,

$$f_{\theta}(\vec{x}) = \sigma_L(W_L \sigma_{L-1}(W_{L-1} \sigma_{L-2}(W_{L-2} \cdots \sigma_1(W_1 \vec{x}) \cdots))), \tag{47}$$

where $W_l \in \mathbb{R}^{n_{l+1},n_l}$, $l=1,2,\cdots,L$ and $\sigma_k(\cdot)$ $(k=1,\cdots)$ is an element-wise non-linear function. We note θ as a set of parameters, namely all elements of W, which determine to minimize a loss function. Note that, including convolutional neural networks, most of known neural networks can be written in this form. It is known that, deep neural networks are universal approximator (4-7), namely which can approximate any maps in desired precision as Fourier expansion.

Let us introduce notation for later purpose. One of ingredients is a linear transformation,

$$z_i^{(l)} = \sum_j w_{ij}^{(l)} u_j^{(l-1)}, \tag{48}$$

and the other is a non-linear function, $u_i^{(l)} = \sigma_l(z_i^{(l)})$. We call $z_i^{(l)}$ as a pre-activation variable. $w_{ij}^{(l)}$ is an element of W_l . $u_i^{(l)}$ is *i*-th component of $\vec{u}^{(l)}$ and $\vec{u}^{(0)} = \vec{x}$. Both transformation is needed to realizes a universal function from a vector space to another vector space (4–7).

To determine the parameters, we use a loss function. A concrete form if the loss function is not necessary but for example for a regression we can take mean square error,

$$L_{\theta}(\mathcal{D}) = \frac{1}{2|\mathcal{D}|} \sum_{(\vec{x}_i, \vec{y}_i) \in \mathcal{D}} (\vec{y}_i - f_{\theta}(\vec{x}_i))^2, \tag{49}$$

where $\vec{y_i}$ is desired answer for $\vec{x_i}$, desired output for $\vec{x_i}$. \mathcal{D} is a set of pairs of data, which is replaced by a part of data for mini-batch training. $|\mathcal{D}|$ is the size of data. This quantified a kind

of distance between probability distribution of the answer and distribution of output. We can choose appropriate a loss function for each problem. Popular choice is the Kullback-Leibler divergence (relative entropy), the cross-entropy and so on (85).

A neural network is not convex function in general, so there are no best way to find an optimal value. Moreover, we cannot know probability distribution which generates data, so we have to use sampling average. Practically, parameters in neural networks are tuned by the gradient descent or its variant,

$$\theta \leftarrow \theta - \eta \frac{\partial L_{\theta}(\mathcal{D})}{\partial \theta},\tag{50}$$

where η is a real positive small number called learning rate and θ represents elements of weight matrices $W^{(l)}$. This is called stochastic gradient descent since it uses sampling approximation to evaluate the value of gradient instead of the exact expectation value. An extreme case, the size of \mathcal{D} is taken to 1, it is called on-line training.

The derivative term $\frac{\partial L_{\theta}(\mathcal{D})}{\partial \theta}$, can be evaluated by a recursive formula called backpropagation and the delta rule.

$$\frac{\partial L_{\theta}(\mathcal{D})}{\partial w_{ij}^{(l)}} = \sum_{k} \frac{\partial L_{\theta}(\mathcal{D})}{\partial z_{k}^{(l)}} \frac{\partial z_{k}^{(l)}}{\partial w_{ij}^{(l)}} = \delta_{i}^{(l)} u_{j}^{(l-1)}, \tag{51}$$

we defined,

$$\delta_i^{(l)} = \frac{\partial L_{\theta}(\mathcal{D})}{\partial z_i^{(l)}}.$$
 (52)

while $z_i^{(l+1)} = \sum_j w_{ij}^{(l+1)} \sigma_l(z_j^{(l)})$, we can relate $\delta^{(l)}$ and $\delta^{(l+1)}$,

$$\delta_i^{(l)} = \sum_j \frac{\partial L_{\theta}(\mathcal{D})}{\partial z_j^{(l+1)}} \frac{\partial z_j^{(l+1)}}{\partial z_i^{(l)}} = \sum_k \delta_k^{(l+1)} \frac{\partial z_k^{(l+1)}}{\partial z_i^{(l)}} = \sum_k \delta_k^{(l+1)} w_{ki}^{(l+1)} \sigma_l'(z_i^{(l)}). \tag{53}$$

These recursive equations enables us to calculate correction very efficiently. The information propagates $\vec{x} = \vec{u}^{(1)} \to \vec{u}^{(2)} \to \cdots \to \vec{u}^{(L)}$. The correction propagate backwards,

$$\vec{\delta}^{(1)} \leftarrow \vec{\delta}^{(2)} \leftarrow \dots \leftarrow \vec{\delta}^{(L)},\tag{54}$$

where $\vec{\delta}^{(l)}$ has $\delta_k^{(l)}$ as elements. This is called backpropagation or backprop in short because error propagates backwards.

A.2 Gauge covariance on the lattice

Here we briefly introduce concept of gauge covariance. A link variable is defined on a bond on a lattice, which is a matrix valued function. A gauge link defined on a point n with a direction μ is transformed like,

$$U_{\mu}(n) \to G(n)U_{\mu}(n)G^{\dagger}(n+\hat{\mu}).$$
 (55)

where G(n) is a SU(N) or U(1) valued function. This transformation law uses different element for each lattice point and it is a local symmetry. In addition, this transformation of gauge link involves two of transformation matrix G(n) and $G^{\dagger}(n+\hat{\mu})$. This is not invariance but rather called covariance.

We can make a Wilson line which obeys same transformation rule with the link. For example,

$$W_{\mu}(n) = U_{\nu}(n)U_{\mu}(n+\hat{\nu})U_{\nu}^{\dagger}(n+\hat{\mu}) \to G(n)W_{\mu}(n)G^{\dagger}(n+\hat{\mu}).$$
 (56)

By using covariant operators, we can construct invariant objects using hermitian conjugate † operation and trace. For example, one non-trivial invariant object a plaquette is given,

$$tr[U_{\mu}(n)W_{\mu}^{\dagger}(n)] \tag{57}$$

and it is transforms like,

$$\operatorname{tr}[U_{\mu}(n)W_{\mu}^{\dagger}(n)] \to \operatorname{tr}[G(n)U_{\mu}(n)G^{\dagger}(n+\hat{\mu})G(n+\hat{\mu})W_{\mu}^{\dagger}(n)G^{\dagger}(n)] = \operatorname{tr}[U_{\mu}(n)W_{\mu}^{\dagger}(n)]$$
(58)

By construction, traced loop operators are invariant because transformation matrices are canceled.

A.3 Tensor calculus and star product

A.3.1 Definition of matrix derivative

The derivative of a scalar f function of a matrix A of independent variables, with respect to the matrix A is a rank-2 tensor:

$$\left[\frac{\partial f[A]}{\partial A}\right]^{i}{}_{j} \equiv \frac{\partial f}{\partial A^{j}{}_{i}}.$$
 (59)

The derivative of a matrix M function of a matrix A of independent variables, with respect to the matrix A is a rank-4 tensor:

$$\left[\frac{\partial M[A]}{\partial A}\right]_{j}^{i}{}_{l}^{k} \equiv \frac{\partial}{\partial A^{j}{}_{k}} M[A]^{i}{}_{l} \tag{60}$$

This is matrix-generalized Wirtinger derivative and for derivative with a complex number, it is fall back to conventional Wirtinger derivative. The derivative of a matrix M function of a scalar x of independent variables, with respect to the matrix x is a rank-2 tensor:

$$\left[\frac{\partial M[x]}{\partial x}\right]^{i}{}_{j} \equiv \frac{\partial M[x]^{i}{}_{j}}{\partial x}.$$
(61)

A.3.2 Chain rules and star products

The chain rule is given as

$$\left[\frac{\partial f[A]}{\partial A}\right]^{i}_{j} = \frac{\partial f}{\partial A^{j}_{i}} = \sum_{kl} \frac{\partial f}{\partial B^{k}_{l}} \frac{\partial B^{k}_{l}}{\partial A^{j}_{i}} \tag{62}$$

$$= \sum_{kl} \left[\frac{\partial f}{\partial B} \right]^{l}{}_{k} \left[\frac{\partial B}{\partial A} \right]^{k}{}_{j}{}^{i}{}_{l} \tag{63}$$

Then, by defining rank-2-rank4 "star"-product

$$[A \star T]^i{}_j \equiv \sum_{kl} A^l{}_k T^k{}_j{}^i{}_l, \tag{64}$$

$$\frac{\partial f}{\partial A} = \frac{\partial f}{\partial B} \star \frac{\partial B}{\partial A} \tag{65}$$

We consider the following chain rule:

$$\left[\frac{\partial M[A]}{\partial A}\right]_{j}^{i}{}_{l}^{k} = \sum_{nm} \frac{\partial M[A]^{i}{}_{l}}{\partial B^{n}{}_{m}} \frac{\partial B^{n}{}_{m}}{\partial A^{j}{}_{k}}$$

$$(66)$$

$$= \sum_{nm} \left[\frac{\partial M[B]}{\partial B}\right]_{n}^{i}{}_{n}{}^{m}{}_{l} \left[\frac{\partial B[A]}{\partial A}\right]_{j}^{n}{}_{m}$$
 (67)

By defining rank-4-rank4 star product:

$$(S \star T)^{i}{}_{j}{}^{k}{}_{l} = \sum_{nm} S^{i}{}_{n}{}^{m}{}_{l} T^{n}{}_{j}{}^{k}{}_{m}$$
 (68)

we have

$$\frac{\partial M[A]}{\partial A} = \frac{\partial M[B]}{\partial B} \star \frac{\partial B[A]}{\partial A} \tag{69}$$

We consider a different chain rule:

$$\frac{\partial f}{\partial x} = \sum_{ij} \frac{\partial f}{\partial A^{j}_{i}} \frac{\partial A^{j}_{i}}{\partial x} \tag{70}$$

$$= \sum_{ij} \left[\frac{\partial f}{\partial A} \right]^{i}_{j} \left[\frac{\partial A}{\partial x} \right]^{j}_{i} \tag{71}$$

$$= \operatorname{Tr}\left[\left[\frac{\partial f}{\partial A}\right]\left[\frac{\partial A}{\partial x}\right]\right]$$
 (72)

Here, f and x are scalars. We can define a rank-2-rank-2 star product.

$$A \star B \equiv \text{Tr} \left[AB \right] \tag{73}$$

And we have

$$\frac{\partial f}{\partial x} = \sum_{ij} \frac{\partial f}{\partial A^{j}_{i}} \frac{\partial A^{j}_{i}}{\partial x} \tag{74}$$

$$= \sum_{ij} \sum_{kl} \frac{\partial f}{\partial B^{k}_{l}} \frac{\partial B^{k}_{l}}{\partial A^{j}_{i}} \frac{\partial A^{j}_{i}}{\partial x}$$
 (75)

$$= \sum_{ij} \sum_{kl} \left[\frac{\partial f}{\partial B} \right]^l{}_k \left[\frac{\partial B}{\partial A} \right]^k{}_i{}^j{}_l \left[\frac{\partial A}{\partial x} \right]^i{}_j \tag{76}$$

$$= \sum_{ij} \sum_{kl} \left[\frac{\partial f}{\partial B} \right]^l{}_k \left[\frac{\partial B}{\partial A} \star \frac{\partial A}{\partial x} \right]^k{}_l \tag{77}$$

$$= \frac{\partial f}{\partial B} \star \left(\frac{\partial B}{\partial A} \star \frac{\partial A}{\partial x} \right) \tag{78}$$

Here we use the following definition:

$$[S \star A]^i{}_j \equiv S^i{}_k{}^l{}_j A^k{}_l \tag{79}$$

This is a definition of a rank-4-rank-2 star product.

A.3.3 Product rules

We consider the following product rule:

$$\left[\frac{\partial (f[A]M[A])}{\partial A}\right]_{j}^{i}{}_{l}^{k} = \frac{\partial}{\partial A^{j}{}_{k}}(f[A]M[A]^{i}{}_{l}]) \tag{80}$$

$$= \frac{\partial f[A]}{\partial A^{j}_{k}} M[A]^{i}_{l} + f[A] \frac{\partial}{\partial A^{j}_{k}} M[A]^{i}_{l}$$
(81)

$$= \left[\frac{\partial f[A]}{\partial A}\right]^{k}{}_{j}M[A]^{i}{}_{l} + f[A]\left[\frac{\partial M[A]}{\partial A}\right]^{i}{}_{j}{}^{k}{}_{l}$$
(82)

By defining "direct" product:

$$(A \oplus B)^{i}{}_{j}{}^{k}{}_{l} = A^{k}{}_{j}B^{i}{}_{l} \tag{83}$$

$$\frac{\partial f[A]M[A]}{\partial A} = \frac{\partial f[A]}{\partial A} \oplus M[A] + f[A]\frac{\partial M[A]}{\partial A}$$
(84)

We consider the matrix M = ABC. The derivative is

$$\left[\frac{\partial M[B]}{\partial B}\right]_{j}^{i}{}_{l}^{k} = \frac{\partial}{\partial B^{j}{}_{k}} M[B]^{i}{}_{l} \tag{85}$$

$$= \sum_{n,m} \frac{\partial}{\partial B^{j}_{k}} (A^{i}_{n} B^{n}_{m} C^{m}_{l}) \tag{86}$$

$$\sum_{n,m} (A^i{}_n \delta_{jn} \delta_{km} C_l^m) \tag{87}$$

$$=A^{i}{}_{j}C^{k}_{l} \tag{88}$$

By defining the outer product

$$(A \otimes B)^{i}{}_{j}{}^{k}{}_{l} \equiv A^{i}{}_{j}B^{k}_{l} \tag{89}$$

we have

$$\frac{\partial M[B]}{\partial B} = A \otimes C \tag{90}$$

We consider

$$\left[\frac{\partial (Y[A]Z[A])}{\partial A}\right]_{j}^{i}{}_{l}^{k}{}_{l} = \frac{\partial}{\partial A^{j}{}_{k}}[YZ]^{i}{}_{l} \tag{91}$$

$$=\sum_{n}\frac{\partial}{\partial A^{j}_{k}}(Y^{i}_{n}Z^{n}_{l})\tag{92}$$

$$= \sum_{n} \frac{\partial Y^{i}_{n}}{\partial A^{j}_{k}} Z^{n}_{l} + \sum_{n} Y^{i}_{n} \frac{\partial Z^{n}_{l}}{\partial A^{j}_{k}}$$

$$\tag{93}$$

$$= \sum_{n} \left[\frac{\partial Y[A]}{\partial A}\right]_{j}^{i}{}_{n}^{k} Z^{n}{}_{l} + \sum_{n} Y^{i}{}_{n} \left[\frac{\partial Z[A]}{\partial A}\right]_{j}^{n}{}_{l}^{k}$$
(94)

By defining the contraction:

$$[AT]_{j}^{i}{}_{l}^{k} = \sum_{n} A_{n}^{i} T_{j}^{n}{}_{l}^{k}$$
(95)

$$[TA]_{j}^{i}{}_{l}^{k} = \sum_{n} T_{j}^{i}{}_{n}A^{n}{}_{l}$$
(96)

$$\frac{\partial YZ}{\partial A} = \frac{\partial Y[A]}{\partial A}Z + Y\frac{\partial Z[A]}{\partial A} \tag{97}$$

A.3.4 Useful formulas

The useful formulas are given as

$$\frac{\partial A}{\partial A} = I \otimes I \tag{98}$$

$$A \star (B \oplus C) = \text{Tr}(AC)B \tag{99}$$

$$A \star (B \otimes C) = CAB \tag{100}$$

$$(A \otimes B) \star (C \otimes D) = (AC) \otimes (DB) \tag{101}$$

$$(A \oplus B) \star (C \oplus D) = \operatorname{Tr}(AD)(C \oplus B) \tag{102}$$

$$(A \otimes B) \star (C \oplus D) = C \oplus (ADB) \tag{103}$$

$$(A \oplus B) \star (C \otimes D) = (DAC) \oplus B \tag{104}$$

$$S \star (T \star K) = (S \star T) \star K \tag{105}$$

$$A \star (T \star K) = (A \star T) \star K \tag{106}$$

$$(A \otimes B)C = A \otimes (BC) \tag{107}$$

$$(A \oplus B)C = A \oplus (BC) \tag{108}$$

$$C(A \otimes B) = (CA) \otimes B \tag{109}$$

$$C(A \oplus B) = A \oplus (CB) \tag{110}$$

$$(S \star T)C = (SC) \star T \tag{111}$$

$$A \star (TC) = (CA) \star T \tag{112}$$

$$S \star (I \otimes I) = S \tag{113}$$

$$B^{\dagger} \star \frac{\partial M^{\dagger}}{\partial A^{\dagger}} = [B \star \frac{\partial M}{\partial A}]^{\dagger} \tag{114}$$

A.4 Backpropagation for stout smearing

A.4.1 Definition

Smearing step in the stout smearing (73) is expressed as

$$U_{\mu}^{(l+1)}(n) = \exp(Q_{\mu}^{(l)}(n))U_{\mu}^{(l)}(n)$$
(115)

$$= U_{\mu}^{(l)}(n) + \left(\exp(Q_{\mu}^{(l)}(n)) - 1\right) U_{\mu}^{(l)}(n). \tag{116}$$

So, we can identify,

$$w_1^{(l)} = w_2^{(l)} = 1 (117)$$

$$\mathcal{N}(z) = z \tag{118}$$

$$\mathcal{G}_{\mu,n}(U^{(l)}) = \left(\exp(Q_{\mu}^{(l)}(n)) - 1\right) U_{\mu}^{(l)}(n) \tag{119}$$

where $Q_{\mu}^{l}(\boldsymbol{n})$ and

$$Q_{\mu}(n) = -\frac{1}{2}(\Omega_{\mu}(n) - \Omega_{\mu}^{\dagger}(n)) + \frac{1}{2N}\operatorname{Tr}(\Omega_{\mu}^{\dagger}(n) - \Omega_{\mu}(n))$$
 (120)

$$\Omega_{\mu}(n) = C_{\mu}(n)U_{\mu}^{\dagger}(n) \tag{121}$$

and $C_{\mu}(n)$ is the weighted sum of the perpendicular staples which begin at lattice site n and terminate at neighboring site $n + \hat{\mu}$. For example, the staple of a plaquette loop is given as

$$\sum_{\nu \neq \nu} \rho_{\mu,\nu} \left(U_{\nu}(n) U_{\mu}(n+\hat{\nu}) U_{\nu}^{\dagger}(n+\hat{\mu}) + U_{\nu}^{\dagger}(n-\hat{\nu}) U_{\mu}(x-\nu) U_{\nu}(x-\hat{\nu}+\hat{\mu}) \right)$$
(122)

A.4.2 Backpropagation

In the stout smearing scheme, the $\delta_{\mu}(n)$ is given as

$$\delta_{\mu}^{(l)}(n) = \frac{\partial S}{\partial z_{\mu}^{(l)}(n)} = \frac{\partial S}{\partial U_{\mu}^{(l)}(n)} \tag{123}$$

$$\bar{\delta}_{\mu}^{(l)}(n) = \frac{\partial S}{\partial U_{\mu}^{(l)\dagger}(n)} = \frac{\partial S}{\partial U_{\mu}^{(l)}(n)} \star \frac{\partial U_{\mu}^{(l)}(n)}{\partial U_{\mu}^{(l)\dagger}(n)} \tag{124}$$

$$= \frac{\partial S}{\partial U_{\mu}^{(l)}(n)} \star \frac{\partial (U_{\mu}^{(l)\dagger}(n))^{-1}}{\partial U_{\mu}^{(l)\dagger}(n)}$$

$$\tag{125}$$

$$= -\frac{\partial S}{\partial U_{\mu}^{(l)}(n)} \star (U_{\mu}^{(l)}(n) \otimes U_{\mu}^{(l)(n)})$$

$$\tag{126}$$

$$= -U_{\mu}^{(l)}(n)\delta_{\mu}^{(l)}(n)U_{\mu}^{(l)}(n) \tag{127}$$

On the l-th layer, the backpropagation is given as

$$\begin{split} & \delta_{\mu}^{(l)}(n) = \delta_{\mu}^{(l+1)}(n) \star (I \otimes I) \\ & + \sum_{\mu',m} \delta_{\mu'}^{(l+1)}(m) \star \left[\frac{\partial \mathcal{G}_{\mu',m}(U^{(l)})}{\partial U_{\mu}^{(l)}(n)} \star (I \otimes I) \right] - \sum_{\mu',m} (U_{\mu',m} \delta_{\mu'}^{(l+1)}(m) U_{\mu',m}) \star \left[\frac{\partial \mathcal{G}_{\mu',m}(U^{(l)})^{\dagger}}{\partial U_{\mu}^{(l)}(n)} \star (I \otimes I) \right] \end{split}$$

By substituting

$$\frac{\partial \mathcal{G}_{\mu',m}(U^{(l)})}{\partial U_{\mu}^{(l)}(n)} = \frac{\partial \exp(Q_{\mu'}^{(l)}(m))}{\partial U_{\mu'}^{(l)}(n)} U_{\mu'}^{(l)}(m) + (\exp(Q_{\mu'}^{(l)}(m)) - I) \frac{\partial U_{\mu'}^{(l)}(m)}{\partial U_{\mu}^{(l)}(n)}$$
(128)

$$= \frac{\partial \exp(Q_{\mu'}^{(l)}(m))}{\partial U_{\mu'}^{(l)}(n)} U_{\mu'}^{(l)}(m) + (\exp(Q_{\mu'}^{(l)}(m)) - I) \delta_{\mu,\mu'} \delta_{n,m} I \otimes I$$
 (129)

$$\frac{\partial \mathcal{G}_{\mu',m}^{\dagger}(U^{(l)})}{\partial U_{\mu}^{(l)}(n)} = U_{\mu'}^{(l)\dagger}(m) \frac{\partial \exp(Q_{\mu'}^{(l)\dagger}(m))}{\partial U_{\mu}^{(l)}(n)}$$
(130)

$$\delta_{\mu}^{(l)}(n) = \delta_{\mu}^{(l+1)}(n)) \exp(Q_{\mu'}^{(l)}(m)) \star (I \otimes I)
+ \sum_{\mu',m} (U_{\mu'}^{(l)}(m)\delta_{\mu'}^{(l+1)}(m)) \star \frac{\partial \exp(Q_{\mu'}^{(l)}(m))}{\partial U_{\mu}^{(l)}(n)} \star (I \otimes I)
- \sum_{\mu',m} (\delta_{\mu'}^{(l+1)\dagger}(m)U_{\mu'}^{(l)\dagger}(m)) \star \left[-\frac{\partial \exp(Q_{\mu'}^{(l)\dagger}(m))}{\partial U_{\mu}^{(l)}(n)} \right] \star (I \otimes I)$$
(131)

Here we use $\delta = -U^{\dagger}\delta^{\dagger}U^{\dagger}$.

The exponentials are given as

$$\frac{\partial \exp(Q_{\mu'}^{(l)}(m))}{\partial U_{\mu}^{(l)}(n)} = \frac{\partial \exp(Q_{\mu'}^{(l)}(m))}{\partial Q_{\mu'}^{(l)}(m)} \star \frac{\partial Q_{\mu'}^{(l)}(m)}{\partial U_{\mu}^{(l)}(n)}
= \frac{\partial \exp(Q_{\mu'}^{(l)}(m))}{\partial Q_{\mu'}^{(l)}(m)} \star \left[\frac{\partial Q_{\mu'}^{(l)}(m)}{\partial \Omega_{\mu}^{(l)}(n)} \star \frac{\partial \Omega_{\mu'}^{(l)}(m)}{\partial U_{\mu}^{(l)}(n)} + \frac{\partial Q_{\mu'}^{(l)}(m)}{\partial \Omega_{\mu}^{(l)\dagger}(n)} \star \frac{\partial \Omega_{\mu'}^{(l)\dagger}(m)}{\partial U_{\mu}^{(l)}(n)} \right]
= \frac{\partial \exp(Q_{\mu'}^{(l)}(m))}{\partial Q_{\mu'}^{(l)}(m)} \star \frac{\partial Q_{\mu'}^{(l)}(m)}{\partial \Omega_{\mu}^{(l)}(n)} \star \left[\frac{\partial \Omega_{\mu'}^{(l)}(m)}{\partial U_{\mu}^{(l)}(n)} - \frac{\partial \Omega_{\mu'}^{(l)\dagger}(m)}{\partial U_{\mu}^{(l)}(n)} \right]$$
(134)

and

$$\frac{\partial \exp(Q_{\mu'}^{(l)\dagger}(m))}{\partial U_{\mu}^{(l)}(n)} = \frac{\partial \exp(Q_{\mu'}^{(l)\dagger}(m))}{\partial Q_{\mu'}^{(l)\dagger}(m)} \star \frac{\partial Q_{\mu'}^{(l)\dagger}(m)}{\partial \Omega_{\mu}^{(l)}(n)} \star \left[\frac{\partial \Omega_{\mu'}^{(l)}(m)}{\partial U_{\mu}^{(l)}(n)} - \frac{\partial \Omega_{\mu'}^{(l)\dagger}(m)}{\partial U_{\mu}^{(l)}(n)} \right]$$
(135)

By using the fact that $Q^{(l)}$ is a traceless anti-hermitian matrix, we have

$$\frac{\partial Q_{\mu'}^{(l)\dagger}(m)}{\partial \Omega_{\mu}^{(l)}(n)} = -\frac{\partial Q_{\mu'}^{(l)}(m)}{\partial \Omega_{\mu}^{(l)}(n)}$$
(136)

Thus, the equations can be expressed as

$$\delta_{\mu}^{(l)}(n) = \delta_{\mu}^{(l+1)}(n) \exp(Q_{\mu'}^{(l)}(m)) \star (I \otimes I)
+ \sum_{\mu',m} (U_{\mu'}^{(l)}(m)\delta_{\mu'}^{(l+1)}(m)) \star \frac{\partial \exp(Q_{\mu'}^{(l)}(m))}{\partial Q_{\mu'}^{(l)}(m)} \star \frac{\partial Q_{\mu'}^{(l)}(m)}{\partial Q_{\mu'}^{(l)}(n)} \star \left[\frac{\partial \Omega_{\mu'}^{(l)}(m)}{\partial U_{\mu}^{(l)}(n)} - \frac{\partial \Omega_{\mu'}^{(l)\dagger}(m)}{\partial U_{\mu}^{(l)}(n)} \right] \star (I \otimes I)
- \sum_{\mu',m} (U_{\mu'}^{(l)}(m)\delta_{\mu'}^{(l+1)}(m))^{\dagger} \star \frac{\partial \exp(Q_{\mu'}^{(l)\dagger}(m))}{\partial Q_{\mu'}^{(l)\dagger}(m)} \star \frac{\partial Q_{\mu'}^{(l)}(m)}{\partial \Omega_{\mu}^{(l)}(n)} \star \left[\frac{\partial \Omega_{\mu'}^{(l)}(m)}{\partial U_{\mu}^{(l)}(n)} - \frac{\partial \Omega_{\mu'}^{(l)\dagger}(m)}{\partial U_{\mu}^{(l)}(n)} \right] \star (I \otimes I)
(137)$$

$$= \delta_{\mu}^{(l+1)}(n) \exp(Q_{\mu'}^{(l)}(m)) \star (I \otimes I)$$

$$+\sum_{\mu',m}(M_{\mu',m}-M_{\mu',m}^{\dagger})\star\frac{\partial Q_{\mu'}^{(l)}(m)}{\partial \Omega_{\mu}^{(l)}(n)}\star\left[\frac{\partial \Omega_{\mu'}^{(l)}(m)}{\partial U_{\mu}^{(l)}(n)}-\frac{\partial \Omega_{\mu'}^{(l)\dagger}(m)}{\partial U_{\mu}^{(l)}(n)}\right]\star(I\otimes I),\tag{138}$$

$$M_{\mu',m} \equiv (U_{\mu'}^{(l)}(m)\delta_{\mu'}^{(l+1)}(m)) \star \frac{\partial \exp(Q_{\mu'}^{(l)}(m))}{\partial Q_{\mu'}^{(l)}(m)}$$
(139)

By using the following equation:

$$\frac{\partial Q_{\mu}(y)}{\partial \Omega_{\mu}(y)} = -\frac{\partial}{\partial \Omega_{\mu}(y)} \left[\frac{1}{2} \left[\Omega_{\mu}(y)^{\dagger} - \Omega_{\mu}(y) \right] - \frac{1}{2N} \text{Tr} \Omega_{\mu}(y)^{\dagger} I + \frac{1}{2N} \text{Tr} \Omega_{\mu}(y) I \right]$$
(140)

$$= \frac{1}{2} \frac{\partial \Omega_{\mu}(y)}{\partial \Omega_{\mu}(y)} - \frac{1}{2N} \frac{\partial \text{Tr}\Omega_{\mu}(y)}{\partial \Omega_{\mu}(y)} \oplus I$$
(141)

$$=\frac{1}{2}I\otimes I - \frac{1}{2N}I \oplus I \tag{142}$$

and the following formula

$$A \star (\frac{1}{2}I \otimes I - \frac{1}{2N}I \oplus I) = \frac{1}{2}A \star (I \otimes I) - \frac{1}{2N}A \star (I \oplus I)$$
 (143)

$$=\frac{1}{2}A - \frac{1}{2N}\operatorname{Tr}(A)I\tag{144}$$

we obtain

$$\delta_{\mu}^{(l)}(n) = \delta_{\mu}^{(l+1)}(n) \exp(Q_{\mu'}^{(l)}(m)) \star (I \otimes I)$$

$$- \left[\partial Q^{(l)}(m) - \partial Q^{(l)\dagger}(m) \right]$$

$$+\sum_{\mu',m} \Lambda_{\mu',m} \star \left[\frac{\partial \Omega_{\mu'}^{(l)}(m)}{\partial U_{\mu}^{(l)}(n)} - \frac{\partial \Omega_{\mu'}^{(l)\dagger}(m)}{\partial U_{\mu}^{(l)}(n)} \right] \star (I \otimes I)$$

$$(145)$$

$$= \delta_{\mu}^{(l+1)}(n) \exp(Q_{\mu'}^{(l)}(m)) + \sum_{\mu',m} \Lambda_{\mu',m} \star \left[\frac{\partial \Omega_{\mu'}^{(l)}(m)}{\partial U_{\mu}^{(l)}(n)} - \frac{\partial \Omega_{\mu'}^{(l)\dagger}(m)}{\partial U_{\mu}^{(l)}(n)} \right]$$
(146)

Here, $\Lambda_{\mu,n}$ is defined as

$$\Lambda_{\mu',m} \equiv \frac{1}{2} (M_{\mu',m} - M_{\mu',m}^{\dagger}) - \frac{1}{2N} \text{Tr}(M_{\mu',m} - M_{\mu',m}^{\dagger}) I$$
 (147)

With the use of the following equations:

$$\frac{\partial \Omega_{\mu'}^{(l)}(m)}{\partial U_{\mu}^{(l)}(n)} = \frac{\partial (C_{\mu',m} U_{\mu'}^{\dagger}(m))}{\partial U_{\mu}^{(l)}(n)} = \frac{\partial C_{\mu',m}}{\partial U_{\mu}^{(l)}(n)} U_{\mu'}^{(l)\dagger}(m)$$
(148)

$$\frac{\partial \Omega_{\mu'}^{(l)\dagger}(m)}{\partial U_{\mu}^{(l)}(n)} = \frac{\partial (U_{\mu'}(m)C_{\mu',m}^{\dagger})}{\partial U_{\mu}^{(l)}(n)} = U_{\mu'}(m)\frac{\partial C_{\mu',m}^{\dagger}}{\partial U_{\mu}^{(l)}(n)} + \frac{\partial U_{\mu'}(m)}{\partial U_{\mu}^{(l)}(n)}C_{\mu',m}^{\dagger}$$
(149)

$$= U_{\mu'}(m) \frac{\partial C_{\mu',m}^{\dagger}}{\partial U_{\mu}^{(l)}(n)} + \delta_{\mu,\mu'} \delta_{n,m} I \otimes I C_{\mu,n}^{\dagger}$$

$$\tag{150}$$

$$= U_{\mu'}(m) \frac{\partial C_{\mu',m}^{\dagger}}{\partial U_{\mu}^{(l)}(n)} + \delta_{\mu,\mu'} \delta_{n,m} I \otimes C_{\mu,n}^{\dagger}$$

$$\tag{151}$$

$$\frac{\partial \Omega_{\mu'}^{(l)}(m)}{\partial U_{\mu}^{(l)}(n)} - \frac{\partial \Omega_{\mu'}^{(l)\dagger}(m)}{\partial U_{\mu}^{(l)}(n)} = \frac{\partial C_{\mu',m}}{\partial U_{\mu}^{(l)}(n)} U_{\mu}^{(l)\dagger}(n) - U_{\mu'}(m) \frac{\partial C_{\mu',m}^{\dagger}}{\partial U_{\mu}^{(l)}(n)} - \delta_{\mu,\mu'}\delta_{n,m}I \otimes C_{\mu,n}^{\dagger} \quad (152)$$

We finally obtain

$$\delta_{\mu}^{(l)}(n) = \delta_{\mu}^{(l+1)}(n) \exp(Q_{\mu'}^{(l)}(m)) - \Lambda_{\mu,n} \star (I \otimes C_{\mu,n}^{\dagger})$$

$$+ \sum_{\mu',m} \Lambda_{\mu',m} \star \left[\frac{\partial C_{\mu',m}}{\partial U_{\mu'}^{(l)}(n)} U_{\mu'}^{(l)\dagger}(m) - U_{\mu'}(m) \frac{\partial C_{\mu',m}^{\dagger}}{\partial U_{\mu}^{(l)}(n)} \right]$$
(153)

$$= \delta_{\mu}^{(l+1)}(n) \exp(Q_{\mu'}^{(l)}(m)) - C_{\mu,n}^{\dagger} \Lambda_{\mu,n}$$

$$+ \sum_{\mu',m} \Lambda_{\mu',m} \star \left[\frac{\partial C_{\mu',m}}{\partial U_{\mu}^{(l)}(n)} U_{\mu'}^{(l)\dagger}(m) - U_{\mu'}(m) \frac{\partial C_{\mu',m}^{\dagger}}{\partial U_{\mu}^{(l)}(n)} \right]$$
(154)

$$= \delta_{\mu}^{(l+1)}(n) \exp(Q_{\mu'}^{(l)}(m)) - C_{\mu,n}^{\dagger} \Lambda_{\mu,n}$$

$$+ \sum_{\mu',m} \left[(U_{\mu'}^{(l)\dagger}(m)\Lambda_{\mu',m}) \star \frac{\partial C_{\mu',m}}{\partial U_{\mu}^{(l)}(n)} - (\Lambda_{\mu',m}U_{\mu'}(m)) \star \frac{\partial C_{\mu',m}^{\dagger}}{\partial U_{\mu}^{(l)}(n)} \right]$$
(155)

$$= \delta_{\mu}^{(l+1)}(n) \exp(Q_{\mu'}^{(l)}(m)) - C_{\mu,n}^{\dagger} \Lambda_{\mu,n}$$

$$+ \sum_{\mu',m} \left[(U_{\mu'}^{(l)\dagger}(m)\Lambda_{\mu',m}) \star \frac{\partial C_{\mu',m}}{\partial U_{\mu}^{(l)}(n)} + (U_{\mu'}^{\dagger}(m)\Lambda_{\mu',m})^{\dagger} \star \frac{\partial C_{\mu',m}^{\dagger}}{\partial U_{\mu}^{(l)}(n)} \right]$$
(156)

The above equation has only one star product. The rank-4 tensors $\frac{\partial C_{\mu',m}}{\partial U_{\mu}^{(l)}(n)}$ and $\frac{\partial C_{\mu',m}}{\partial U_{\mu}^{(l)}(n)}$ are usually

expressed as

$$\frac{\partial C_{\mu',m}}{\partial U_{\mu}^{(l)}(n)} \equiv \sum_{i} A_{\mu,m}^{i} \otimes B_{\mu,m}^{i} \tag{157}$$

$$\frac{\partial C^{\dagger}_{\mu',m}}{\partial U^{(l)}_{\mu}(n)} \equiv \sum_{i} \bar{A}^{i}_{\mu,m} \otimes \bar{B}^{i}_{\mu,m} \tag{158}$$

Then, we have

$$\delta_{\mu}^{(l)}(n) = \delta_{\mu}^{(l+1)}(n) \exp(Q_{\mu'}^{(l)}(m)) - C_{\mu,n}^{\dagger} \Lambda_{\mu,n} + \sum_{\mu',m} \sum_{i} \left[B_{\mu,m}^{i} U_{\mu'}^{(l)\dagger}(m) \Lambda_{\mu',m} A_{\mu,m}^{i} - \bar{B}_{\mu,m}^{i} \Lambda_{\mu',m} U_{\mu'}(m) \bar{A}_{\mu,m}^{i} \right]$$
(159)

This equation is equivalent to the original one derived without tensor calculations.

A.4.3 parameter derivative

The derivative with respect to the parameter ρ can be calculated as

$$\frac{\partial S}{\partial \rho^{(l)}} = \sum_{\mu,n} \left[\frac{\partial S}{\partial z_{\mu}^{(l)}(n)} \star \frac{\partial z_{\mu}^{(l)}(n)}{\partial \rho^{(l)}} + \frac{\partial S}{\partial z_{\mu}^{(l)\dagger}(n)} \star \frac{\partial z_{\mu}^{(l)\dagger}(n)}{\partial \rho^{(l)}} \right]$$
(160)

$$= \sum_{\mu,n} \left[\delta_{\mu}(n) \star \left(\frac{\partial \exp Q_{\mu}^{(l)}(n)}{\partial \rho^{(l)}} U_{\mu}(n) \right) + \bar{\delta}_{\mu}(n) \star \left(U_{\mu}^{\dagger}(n) \frac{\partial \exp Q_{\mu}^{(l)\dagger}(n)}{\partial \rho^{(l)}} \right) \right]$$
(161)

$$= \sum_{\mu,n} \left[(U_{\mu}^{(l)\dagger}(n)\Lambda_{\mu,n}) \star \frac{\partial C_{\mu,n}}{\partial \rho^{(l)}} + (U_{\mu}^{\dagger}(n)\Lambda_{\mu,n})^{\dagger} \star \frac{\partial C_{\mu,n}^{\dagger}}{\partial \rho^{(l)}} \right]$$
(162)

$$=2\sum_{\mu,n}\operatorname{Re}\operatorname{Tr}\left[U_{\mu}^{(l)\dagger}(n)\Lambda_{\mu,n}\frac{\partial C_{\mu,n}}{\partial\rho^{(l)}}\right]$$
(163)

Here, we use $\frac{\partial C_{\mu,n}^{\dagger}}{\partial
ho^{(l)}} = [\frac{\partial C_{\mu,n}}{\partial
ho^{(l)}}]^{\dagger}$.