# Criticality analysis of nuclear binding energy neural networks

S. A. Sundberg[*] and R. J. Furnstahl[†]

*Department of Physics, The Ohio State University, Columbus, OH 43210, USA*

(Dated: August 5, 2025)

Machine learning methods, in particular deep learning methods such as artificial neural networks (ANNs) with many layers, have become widespread and useful tools in nuclear physics. However, these ANNs are typically treated as "black boxes", with their architecture (width, depth, and weight/bias initialization) and the training algorithm and parameters chosen empirically by optimizing learning based on limited exploration. We test a non-empirical approach to understanding and optimizing nuclear physics ANNs by adapting a criticality analysis based on renormalization group flows in terms of the hyperparameters for weight/bias initialization, training rates, and the ratio of depth to width. This treatment utilizes the statistical properties of neural network initialization to find a generating functional for network outputs at any layer, allowing for a path integral formulation of the ANN outputs as a Euclidean statistical field theory. We use a prototypical example to test the applicability of this approach: a simple ANN for nuclear binding energies. We find that with training using a stochastic gradient descent optimizer, the predicted criticality behavior is realized, and optimal performance is found with critical tuning. However, the use of an adaptive learning algorithm leads to somewhat superior results without concern for tuning and thus obscures the analysis. Nevertheless, the criticality analysis offers a way to look within the black box of ANNs, which is a first step towards potential improvements in network performance beyond using adaptive optimizers.

## I. INTRODUCTION

Machine learning methods such as artificial neural networks (ANNs) offer nuclear physicists new ways to explore and understand nuclear systems, as well as improve existing solution methods [1–3]. Examples of ANN applications include the extrapolation of no-core shell model (NCSM) observables from smaller model space calculations [4], variational Monte Carlo with neural network quantum states [5], and learning nuclear mass models from data [6]. ANNs are commonly treated as black boxes that are empirically optimized. Given their growing prominence as a tool for nuclear physics research, it is desirable to have a framework that allows for a more structured analysis based on an understanding of how they work. Several authors have proposed a field theory approach to analyze and optimize neural networks [7–16], using a combination of methods from quantum field theory and Bayesian statistics. Here we assess whether this approach works in practice for the test case of a feed-forward ANN trained to predict binding energies in the nuclear landscape.

An ANN uses connected neurons arranged in layers that, when trained, give a functional relationship between the input data vector $x$ and the known output data vector $y$, which are entries in a data set $\mathcal{D}$ [17–20]. We use the notation $x_\alpha$ and $y_\alpha$ to refer to a particular entry $\alpha$ in $\mathcal{D}$, while $x_{i,\alpha}$ and $y_{i,\alpha}$ are components of these vectors. In Fig. 1 we show a schematic of a feed-forward ANN with two inputs (so $i$ in $x_{i,\alpha}$ is 1 or 2) and one output (so $i$ in $y_{i,\alpha}$ is 1). The outputs for a neural
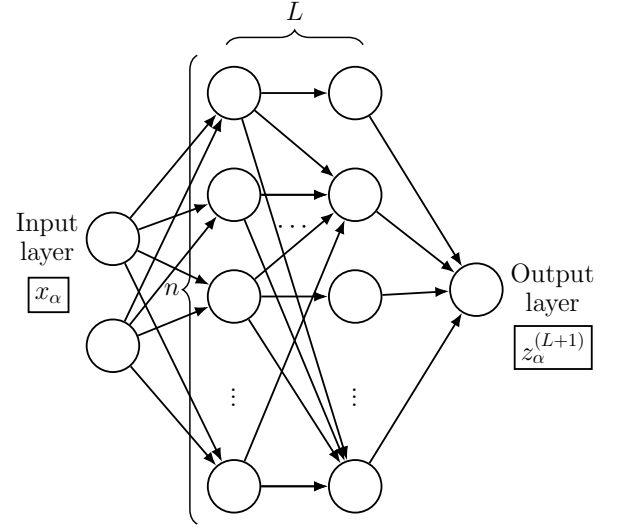


FIG. 1. Schematic of a feed-forward neural network with two inputs $x_\alpha$ and one output $z_\alpha^{(L+1)}$, where $\alpha$ is an index running over the dataset $\mathcal{D}$. The width is $n$ and the depth is $\ell_{out} = L+1$, with $L$ being the number of hidden layers. The neurons are totally connected (some lines are omitted here for clarity).

network are determined by the weights $W_{ij}$ and biases $b_i$, denoted collectively by the parameters $\theta$. On a more "macroscopic level", outputs are determined by the hyperparameters of the network such as the depth $\ell_{out}$, the width $n$, the activation function $\sigma$ (these three hyperparameters are defined here as the architecture hyperparameters) and the initialization hyperparameters that determine the pre-training distributions for the weights and biases. To simplify the analysis, weights and biases
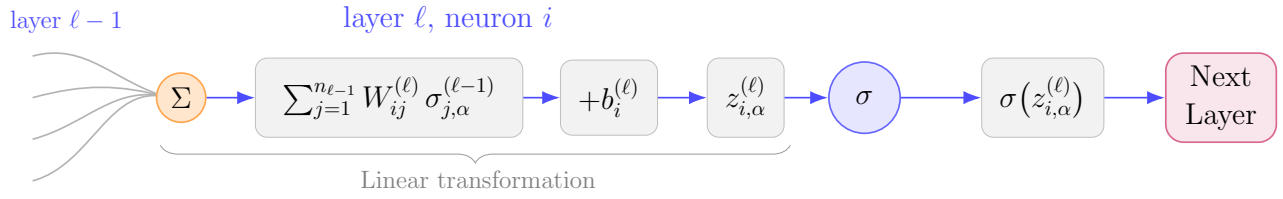
---

[*] sundberg.24@osu.edu
[†] furnstahl.1@osu.edu

FIG. 2. Schematic of a neuron within a neural network. In the layer $\ell$, the product between the weight matrix $W^{(\ell)}$ and the previous layers' output $\sigma^{(\ell-1)}$ is summed over the width of the prior layer $n_{\ell-1}$, and a bias vector $b^{(\ell)}$ is added, forming the preactivation $z^{(\ell)}$. The activation function $\sigma$ acts on the preactivation to yield the output $\sigma(z^{(\ell)})$ (also notated $\sigma^{(\ell)}$) to be passed to the next layer.

will be initialized according to mean-zero Gaussians, so only the variances will be used as initialization hyperparameters. In addition, the width $n_\ell$ in a layer $\ell$ will be uniform across all hidden layers; the length of $x$ will be $n_0$ and the length of $y$ will be $n_{\ell_{out}}(\ell_{out} \equiv L + 1)$.

In Fig. 2, the process by which layer outputs are passed between layers is outlined. Going from the input layer to the first hidden layer ($\ell = 1$) we define $\sigma_{i,\alpha}^{(0)} \equiv x_{i,\alpha}$. For all subsequent layers, a linear transformation maps the outputs from layer $\ell-1$ into the preactivations $z_{i,\alpha}^{(\ell)}$ in layer $\ell$ followed by a transformation by the activation function $\sigma$. For the final hidden layer, the output $z_\alpha^{(\ell_{out})} \equiv f(x_\alpha, \theta)$ is produced with only a linear transformation.

For a given initialization of the network parameters $\theta$, the preactivations and the output function $f$ are uniquely determined. However, repeated random initializations of $\theta$ induces a probability distribution on preactivations and $f$. Explicit expressions for these probability distributions in the limit of large width $n$ (from Refs. [7–10]) will be reviewed in Sec. II, with the ratio of depth to width revealed as an expansion parameter for the distribution's action. The output distribution in a given layer $\ell$ refers to the distribution of preactivations $z_{j,\alpha}^{(\ell)}$. This choice allows for parallel treatment of all other layers with the distribution of functions $f$ in the output layer.

To achieve the desired functional relationship between $x_{i,\alpha}$ and $y_{j,\alpha}$, a neural network is trained by adjusting the weights and biases in each layer according to the gradients of the loss function $\mathcal{L}(f(x, \theta), y)$. The loss function quantifies the difference between the network output $f(x, \theta)$ and the output data $y$, and the goal of training is to minimize $\mathcal{L}(f(x, \theta), y)$. The gradients of $\mathcal{L}(f(x, \theta), y)$ with respect to the network parameters $\theta$ depend on the size of the output, as well as the architecture hyperparameters. It is known empirically that poor choices of architecture and/or initialization lead to exploding/vanishing gradients in the absence of adaptive learning algorithms, and this was problematic for training networks for many years [8, 19, 21].

In modern practice, standard adaptive learning algorithms, plus experience in which activation functions and initializations offer good network training performance, have led routinely to satisfactory results for ANNs us-ing standard libraries. This leaves experimentation with widths and depths of networks as the remaining knobs to adjust for many practitioners [21–23]. However, a more general theory of network behavior offers a way to determine what works and does not work ab initio, to see what network behaviors arise from which hyperparameter, and to identify what improvements can be made. We follow the lead of Refs. [7–10], who construct a theory of neural networks by considering the statistics of the network, calculating a generating functional distribution for moments of the network's output distribution, and analyzing this distribution in analogy to quantum field theory. While this is highly promising, it is also largely theoretical and we seek to test whether it will be of practical use for problems in nuclear physics.

The plan for this paper is as follows. The properties and methods of artificial neural network field theory will be reviewed in Sec. II, then central claims on initialization and learning are validated in Sec. III and Sec. IV, respectively. As a test case, we examine two-input ANNs trained using AME2020 binding energy data. Our conclusions and plans for future studies are summarized in Sec. V.

## II. BASICS OF ANNFT

### A. Overview

The goal of artificial neural network field theory (ANNFT) is to analyze networks using the tools of effective field theory and the renormalization group. This entails working with degrees of freedom based on the output probability distributions of the network in a controlled expansion in the ratio of the depth to width of the ANN.

The ANN output function $f(x, \theta)$ for any fixed $\theta$ (that is, specified weights and biases) is a deterministic function. But, as already stressed, with repeated initializations of the ANN with $\theta$ drawn from random distributions, the network acquires a probability distribution over functions $p(f(x) | \mathcal{D})$. Again, $\mathcal{D}$ is used to represent the data set of inputs and outputs that the network is finding a relationship between. In the pre-training case, this

$$\underbrace{\int d\theta\, P(\theta)\, e^{\int d^d x\, J(x) f(x,\theta)}}_{\text{parameter space}} \implies Z[J] = \left\langle e^{\int d^d x\, J(x) f(x)} \right\rangle \impliedby \underbrace{\int \mathcal{D}f\, e^{-S[f]+\int d^d x\, J(x) f(x)}}_{\text{function space}}$$

FIG. 3. An ensemble of ANNs and a path integral in field theory are both distributions over random functions. The partition function $Z[J]$ is a generating functional for the moments of functions $f$.

quantity simply refers to the inputs $x_\alpha$ for the network. The outputs $y_\alpha$ will become relevant later during training. As we illustrate in Secs. II B and II C, the output probability distribution in any layer for any initialization can be calculated by marginalizing over the weights and biases of the network. The network's dependence on the weights and biases is then transformed into a dependence on the preactivations $z^{(\ell)}$.

In this context, the weights and biases are independently and identically distributed (i.i.d.) random variables, which means the sum defining preactivations $z_{i,\alpha}^{(\ell)}$ in Fig. 2 as the hidden layer width $n$ is taken to be large satisfies the conditions for the central limit theorem (CLT) (assuming mean-zero distributions with appropriate normalization of the variances to keep the sum finite). This implies that the preactivations in each layer and the function output will tend toward a (correlated) Gaussian distribution. In short, the distribution over functions becomes a Gaussian process (GP) in the $n \to \infty$ limit [24–27]. Note that the activation functions do not disturb this conclusion because a nonlinear function of an i.i.d. random variable is still i.i.d.

At a large but finite width, the distribution over output functions acquires non-Gaussian components, but the near-CLT averaging enables these correlations to be treated in a controlled manner. At the same time, increasing depth $\ell_{out}$ will be seen to introduce increasing correlations. These opposing tendencies lead to the identification of the ratio $r \equiv \ell_{out}/n$ being significant as both $n$ and $\ell_{out}$ become large. $r$ acts as an expansion parameter that suppresses higher-order, non-Gaussian correlations in the limit of small $r$. The nature of width and depth, their ratio $r$ as an expansion parameter, and perturbative non-Gaussianity of distributions will be explored Secs. II B, II C, II D, and III.

A complete description of the probability distributions in ANNFT can be formulated in terms of their moments, which characterize the shape of a distribution. The $n^{\text{th}}$-order moment of a distribution of the (untrained) function $f$, $\langle f(x_1) \cdots f(x_n) \rangle$, is the expectation value of this product of functions at specified inputs $x_i$. Examples of moments include the mean, variance, skewness, and kurtosis of a univariate probability distribution (as well as their multivariate distribution counterparts), which are, respectively, the first, second, third, and fourth moments of the distribution. The partition function $Z[J]$ defined in the center of Fig. 3 is the generating function for the

moments, for example,

$$\langle f(x_1) f(x_2) \rangle = \left. \frac{\delta}{\delta J(x_1)} \frac{\delta}{\delta J(x_2)} Z[J] \right|_{J=0}, \qquad (1)$$

so $Z[J]$ is a fundamental object to calculate. In the limit of infinite $n$ (when the distribution becomes Gaussian) the only independent moments are the mean (which will be zero before training if the initialization distributions are mean zero, along with all other odd-ordered moments) and the covariance. In particular, moments of order greater than two are dominantly given by factorized products of the covariance, whereas non-Gaussian contributions will be suppressed by powers of $1/n$. We will take advantage of this simplification and expand about this limit.

Because the distribution of $f$ is determined by the distributions of the weights and biases in $\theta$, one way to calculate the partition function is to evaluate the integral over this parameter space as on the left of Fig. 3, where $P(\theta)$ is the joint probability of all the weights and biases. We will review how to formulate and evaluate this integral by considering successive layers of the ANN in Secs. II B, II C, and II D. The formulation of $Z[J]$ in these parameter-space degrees of freedom (dofs) would seem at first to suggest that the complexity of the network should become very high in the large width and depth limit. However, one finds that the smoothing effect from the approach to the CLT implies that this is not the correct characterization of complexity for the ANN [8].

An alternative is to express the statistical properties of ensembles of neural networks in a path-integral-like formulation, where the integration is over the distribution of random functions (see the right equation in Fig. 3). If applied to a given layer, this generating functional can be used to calculate any correlation function for that layer; if applied to the output layer, this yields correlation functions in terms of the network outputs. By choosing the initialization distributions to be mean-zero Gaussians, the probability distribution for the preactivations at any layer $\ell$ (including the output) takes the form of an exponentiated action:

$$p\big(z^{(\ell)} \,|\, \mathcal{D}\big) \propto \exp\Big(-S(z^{(\ell)})\Big), \qquad (2)$$

where we have omitted the normalization constant. Derivations of the moments of the output function (and more generally the preactivations for each layer) can be simplified with Wick's theorem and diagrammatics. Observables in a given layer can be calculated just like in

quantum field theory with perturbation theory. This depends on having a small ratio of depth to width $r$, and using that as an expansion parameter for the network (see Chapters 0 and $\epsilon$ in Ref [8] for a broader overview).

In the $n \to \infty$ limit, the non-Gaussian contributions are suppressed, and we get a Gaussian free field theory (i.e., a quadratic action). By backing off this limit, the output distribution in a given layer becomes nearly Gaussian. This introduces connected (non-Gaussian) contributions from moments higher than second order, which can be treated perturbatively. The diagrams mentioned before can now include vertices representing the connected contributions, and observables for any layer can be calculated. All correlations of interest are included in an action $S(z)$, each with data-dependent couplings. These couplings now provide a measure of complexity for the network.

The couplings will be shown in Sec. II D to evolve under recursion relations with depth, which act like a renormalization group flow. Increasing depth introduces correlations via this RG flow, while increasing width reduces them through CLT averaging. As such, the terms in the action will use the ratio of depth to width $r = \ell_{out}/n$ as an expansion parameter, controlling which couplings contribute meaningfully to the output distribution of a network. By analyzing these distributions, values of the initialization widths $C_W$ and $C_b$ (defined in Sec. II B) can be found that make the variance constant under the RG flow with depth and lead to higher-order interactions with increasing powers of $r$. These values of the initialization hyperparameters are known as critical values, and by tuning the initialization to criticality, the network output magnitude will not explode or vanish with depth. This in turn prevents the gradients of the network from exploding or vanishing, providing stability in training through the criticality analysis of the networks with ANNFT.

So far we have only considered the initialization of the neural network, with no training with respect to the $y$'s in $\mathcal{D}$. In carrying out such training, the original parameters are adjusted to a new set of parameters $\theta^*$ that minimize the loss function $\mathcal{L}(f(x, \theta), y)$ for each initialization of $\theta$. This implies a distribution $p(f(x, \theta^*) \,|\, \mathcal{D})$ for trained functions $f(x, \theta^*)$, which provides an ensemble of solutions for a given problem and enables uncertainty quantification.

The strategy will be to Taylor expand about initialization with large (but finite) width. If the trained parameters $\theta^*$ are close to the initialized values $\theta$, then we can truncate

$$f(x; \theta^\star) = f(x; \theta) + (\theta^\star - \theta)\frac{df}{d\theta} + \frac{1}{2}(\theta^\star - \theta)^2 \frac{d^2 f}{d\theta^2} + \cdots \tag{3}$$

to good approximation. This gives a way to understand solutions in terms of the initialization, but if the Taylor series has a large number of relevant terms dependent on derivatives of $f(x, \theta)$ to high order, the high dimensionality of the parameters makes such an analysis intractable.

However, Ref. [8] shows that the Taylor series only needs up to the third derivative to describe network outputs with large $n$. The higher derivatives are suppressed by powers of $r = \ell_{out}/n$ in this limit. With these higher-order terms suppressed, the critical initialization helps to ensure that the covariance of the output in any layer does not exponentially grow or vanish with depth, and therefore the magnitudes of the loss function gradients used to update $\theta$ are not too large or small, as they depend on the magnitude of $z^{(l)}$. By regulating the magnitude of the gradients and, therefore, the difference between the initial and trained parameters, critical initialization partially ensures the validity of the Taylor series.

To fully ensure the validity, critical training must be implemented to provide additional regulation of the gradients and the size of the trained parameters. The dependence of the Taylor series on the derivatives, as well as the information about the trained parameters, can be folded into what is called the neural tangent kernel (NTK) of a given network [8, 28, 29]. The NTK is a matrix that governs how the outputs in each layer evolve during training, and will be discussed in more detail in Sec. II E. Like the couplings of the output distribution, the NTK also evolves under an RG flow with depth. Training is tuned to criticality by scaling the learning rates of the network so that the NTK evolves in a controllable, non-exponential fashion like the moments of the distribution. By having the trained output expressed in terms of its dependence on the statistical moments of the initialized output distribution and the NTK, the Taylor series is cast into a form dependent on $z^{(L+1)}(T)$ (the final layer output after $T$ training epochs) and the NTK $H^{(\ell)}$ rather than $f(x, \theta)$ and its many derivatives. The Taylor series can then be truncated by the dependence of the stochastic variables $z^{(\ell)}$ and $H^{(\ell)}$ on the width $n$ and depth $\ell_{out}$.

Significantly, the output distribution after training also becomes a Gaussian distribution in the infinite-width limit, and nearly Gaussian with connected correlations that depend on $r$ in the large-but-finite limit. It can be shown in the limit of small $r$ that the Taylor series can be truncated to a few terms due to this dependence, allowing for predictions of trained network behavior through tractable perturbative calculations.

In the following sections, we fill in selected details of this overview. For much more extensive treatments, the reader is directed to Refs. [8–10].

## B. First hidden layer

To derive the statistical field theory for network outputs, we first consider the statistics of the network parameters. A network's weights $W_{ij}$ and biases $b_i$ are initialized from probability distributions. For the sake of simpler calculations, these are chosen to be mean-zero

Gaussian distributions for the weights and biases:[1]

$$p(W_{ij}^{(\ell)} \mid I) = \sqrt{\frac{n}{2\pi C_W}} \exp\left(-\frac{n(W_{ij}^{(\ell)})^2}{2C_W}\right), \qquad (4)$$

$$p(b_i^{(\ell)} \mid I) = \sqrt{\frac{1}{2\pi C_b}} \exp\left(-\frac{(b_i^{(\ell)})^2}{2C_b}\right), \qquad (5)$$

where $n$ is the common width of the neural network hidden layers (previously the width $n_\ell$ was used to describe the width in a layer $\ell$), and the quantities $C_W$ and $C_b$ are the variances of the distributions. To keep the formulas for parameter distributions compact, we use "$I$" (for information) to represent these contingent quantities in probability density functions (pdfs). Note that while the true variance of the weight distribution is $C_W/n$, the unscaled variance $C_W$ and the bias variance $C_b$ will be the more relevant quantities for calculation and understanding network behavior. The factor of $1/n$ serves to cancel the $n$ terms summed in the previous layer, so that we will have a useful large $n$ limit (as can be seen from the expressions for the preactivations in Fig. 2 or Eq. (6)). Correspondingly, the variance for the weights between the input layer and the first hidden layer is scaled by a factor $1/n_0$ rather than $1/n$.

For a neural network of a given architecture, the probability distribution of preactivations $z_{i,\alpha}^{(\ell)}$ at a given layer $\ell$ from the initialization distributions (4) and (5) can be calculated directly. The preactivations $z^{(\ell)}$ at layer $\ell$ are

(see Fig. 2)

$$z_{i,\alpha}^{(\ell)} = \sum_{j=1}^{n_{\ell-1}} W_{ij}^{(\ell)} \sigma_{j,\alpha}^{(\ell-1)} + b_i^{(\ell)}, \qquad (6)$$

where

$$\sigma_{i,\alpha}^{(\ell)} = \sigma(z_{i,\alpha}^{(\ell)}), \qquad (7)$$

and

$$\sigma_{i,\alpha}^{(0)} = x_{i,\alpha}. \qquad (8)$$

In the first hidden layer ($\ell = 1$), the joint probability distribution $p(z^{(1)} \mid \mathcal{D})$ can be calculated by integrating over the model parameters from $-\infty$ to $\infty$,[2]

$$p(z^{(1)} \mid \mathcal{D}) = \int_{-\infty}^{\infty} \left[\prod_i db_i^{(1)} p(b_i^{(1)})\right] \left[\prod_{i,j} dW_{ij}^{(1)} p(W_{ij}^{(1)})\right]$$
$$\times \prod_{i,\alpha} \delta\left(z_{i;\alpha}^{(1)} - b_i^{(1)} - \sum_j W_{ij}^{(1)} x_{j;\alpha}\right). \quad (9)$$

A strategy for evaluating this integral is to use the integral representation of the delta function[3]

$$\delta(z - a) = \int_{-\infty}^{\infty} \frac{d\Lambda}{2\pi} e^{i\Lambda(z-a)}. \qquad (10)$$

Applying this to all indices $\delta\left(z_{i;\alpha} - b_i - \sum_j W_{ij} x_{j;\alpha}\right)$ results in Eq. (9) becoming

$$p(z^{(1)} \mid \mathcal{D}) = \int_{-\infty}^{\infty} \left[\prod_i \frac{db_i^{(1)}}{\sqrt{2\pi C_b}}\right] \left[\prod_{i,j} \frac{dW_{ij}^{(1)}}{\sqrt{2\pi C_W/n_0}}\right] \left[\prod_{i,\alpha} \frac{d\Lambda_i^\alpha}{2\pi}\right]$$
$$\times \exp\left[-\sum_i \frac{(b_i^{(1)})^2}{2C_b} - n_0 \sum_{i,j} \frac{(W_{ij}^{(1)})^2}{2C_W} + i\sum_{i,\alpha} \Lambda_i^\alpha \left(z_{i;\alpha} - b_i^{(1)} - \sum_j W_{ij}^{(1)} x_{j;\alpha}\right)\right]. \qquad (11)$$

Evaluating the integrals by completing the square for $W_{ij}$ and $b_i$ yields the distribution

$$p(z^{(1)} \mid \mathcal{D}) = \frac{1}{|2\pi G^{(1)}|^{\frac{n_1}{2}}} e^{-\frac{1}{2} \sum_{i=1}^{n_1} \sum_{\alpha_1,\alpha_2 \in \mathcal{D}} z_{i;\alpha_1}^{(1)} G_{(1)}^{\alpha_1 \alpha_2} z_{i;\alpha_2}^{(1)}}, \qquad (12)$$

where $G_{\alpha_1,\alpha_2}^{(1)}$ is the covariance matrix for the distribution, defined by

$$G_{\alpha_1 \alpha_2}^{(1)} = C_b + C_W \sum_j \frac{x_{j;\alpha_1} x_{j;\alpha_2}}{n_0}, \qquad (13)$$

---

[1] The initialization distributions need not be mean zero, or even Gaussian at all. This choice makes calculations easier to perform, and enhances the analogy to field theories in physics later.

[2] Equation (9) follows because the weights and biases are independent and $p(z_{i,\alpha}^{(\ell)} \mid b_i^{(\ell)}, W_{i,j}^{(\ell)})$ is a delta function for any $\ell$.

[3] The introduction of $\Lambda$ is the analog of the Hubbard-Stratonovich method used in quantum field theory [8].

and the notation $G_{(1)}^{\alpha_1\alpha_2}$ with raised indices denotes the inverse of the covariance matrix.

## C. Second hidden layer

In the first layer, all statistical correlations between neurons are governed entirely by the covariance matrix, making it a multivariate Gaussian distribution. Distributions in further layers develop non-Gaussianity for finite $n$. To see this explicitly, the second layer output distribution $p(z^{(2)}|\mathcal{D})$ can be calculated, and the method of calculation can be extrapolated for all deeper layers. By application of the product rule, and then the sum rule for probability distributions, $z^{(1)}$ can be marginalized over. From this, $p(z^{(2)}|\mathcal{D})$ can be expressed as

$$p(z^{(2)}|\mathcal{D}) = \int_{-\infty}^{\infty} \left[\prod_{j,\alpha} dz_{j;\alpha}^{(1)}\right] p(z^{(2)}|z^{(1)}) p(z^{(1)}|\mathcal{D}),$$
(14)

with the conditional distribution

$$p(z^{(2)}|z^{(1)}) = \int_{-\infty}^{\infty} \left[\prod_i db_i^{(2)} p(b_i^{(2)})\right]\left[\prod_{i,j} dW_{ij}^{(2)} p(W_{ij}^{(2)})\right]$$
$$\times \prod_{i,\alpha} \delta\left(z_{i;\alpha}^{(2)} - b_i^{(2)} - \sum_j W_{ij}^{(2)} \sigma_{j;\alpha}^{(1)}\right).$$
(15)

The conditional distribution can be evaluated the same way as Eq. (9), resulting in the equation

$$p(z^{(2)}|z^{(1)}) = \frac{1}{|2\pi\widehat{G}^{(2)}|^{\frac{n_2}{2}}} e^{-\frac{1}{2}\sum_{i=1}^{n_2}\sum_{\alpha_1,\alpha_2\in\mathcal{D}} z_{i;\alpha_1}^{(2)} \widehat{G}_{(2)}^{\alpha_1\alpha_2} z_{i;\alpha_2}^{(2)}},$$
(16)

with the covariance matrix

$$\widehat{G}_{\alpha_1\alpha_2}^{(2)} \equiv C_b^{(2)} + C_W^{(2)} \frac{1}{n_1}\sum_{j=1}^{n_1} \sigma_{j;\alpha_1}^{(1)} \sigma_{j;\alpha_2}^{(1)}.$$
(17)

Significantly, the covariance matrix in the second layer conditional distribution and all deeper layer conditional distributions is stochastic, as the activation function $\sigma^{(1)} \equiv \sigma^{(1)}(z^{(1)})$ is a function of the random variable $z^{(1)}$ from the previous layer. The use of a hat on a quantity, such as $\widehat{G}_{\alpha_1\alpha_2}^{(2)}$, indicates that it is a function of random variables. The mean of the conditional distribution's covariance $\widehat{G}_{\alpha_1\alpha_2}^{(2)}$ is

$$G_{\alpha_1\alpha_2}^{(2)} \equiv C_b^{(2)} + C_W^{(2)} \langle \sigma_{\alpha_1}\sigma_{\alpha_2}\rangle_{G^{(1)}}.$$
(18)

$\langle g(z)\rangle_G$ indicates a Gaussian expectation value of the function $g(z)$ for a mean zero Gaussian distribution of $z$

with a covariance matrix $G$. As such, $G_{\alpha_1\alpha_2}^{(2)}$ in Eq. (18) does not have a hat as it is the mean of the stochastic covariance $\widehat{G}_{\alpha_1\alpha_2}^{(2)}$. Due to the stochastic nature of the covariance, there will be fluctuations in value around $G_{\alpha_1\alpha_2}^{(2)}$. These fluctuations can be expressed as

$$\Delta\widehat{G}_{\alpha_1\alpha_2}^{(2)} \equiv \widehat{G}_{\alpha_1\alpha_2}^{(2)} - G_{\alpha_1\alpha_2}^{(2)}$$
$$= C_W^{(2)} \frac{1}{n_1}\sum_{j=1}^{n_1} \left(\sigma_{j;\alpha_1}^{(1)} \sigma_{j;\alpha_2}^{(1)} - \langle\sigma_{\alpha_1}\sigma_{\alpha_2}\rangle_{G^{(1)}}\right).$$
(19)

Calculating the mean fluctuations of the covariance $\left\langle \Delta\widehat{G}_{\alpha_1\alpha_2}^{(2)}\right\rangle_{G^{(1)}}$ results in the first term in Eq. (19) to be equal to the second, so that the mean fluctuations of the conditional distribution's covariance are zero. However, the variance of these fluctuations is non-zero,

$$\langle\Delta\widehat{G}_{\alpha_1\alpha_2}^{(2)}\Delta\widehat{G}_{\alpha_3\alpha_4}^{(2)}\rangle_{G^{(1)}} = \frac{1}{n_1}\left(C_W^{(2)}\right)^2\left[\langle\sigma_{\alpha_1}\sigma_{\alpha_2}\sigma_{\alpha_3}\sigma_{\alpha_4}\rangle_{G^{(1)}}\right.$$
$$\left. - \langle\sigma_{\alpha_1}\sigma_{\alpha_2}\rangle_{G^{(1)}}\langle\sigma_{\alpha_3}\sigma_{\alpha_4}\rangle_{G^{(1)}}\right]$$
$$\equiv \frac{1}{n_1}V_{(\alpha_1\alpha_2)(\alpha_3\alpha_4)}^{(2)}.$$
(20)

The quantity $V_{(\alpha_1\alpha_2)(\alpha_3\alpha_4)}^{(2)}$ is the non-Gaussian, or connected contribution to the fourth-order moment, and in analogy with field theory it serves as the four-point vertex for neuron interactions. As can be seen, the fluctuations in the covariance matrix of a Gaussian distribution depend on the width of the prior layer $n_1$. In the large-width limit, $\frac{1}{n_1}V_{(\alpha_1\alpha_2)(\alpha_3\alpha_4)}^{(2)}$ is suppressed. This serves as a clear example of simplification in the limit of large width, and allows for the vertex $V_{(\alpha_1\alpha_2)(\alpha_3\alpha_4)}^{(2)}$ to be treated perturbatively.

At infinite width, the fluctuations of $p(z^{(2)}|z^{(1)})$'s covariance matrix vanish, and it becomes a deterministic Gaussian distribution. With the simplifications from large $n$, the full second-layer distribution $p(z^{(2)}|\mathcal{D})$ can be calculated. The $\ell = 2$ stochastic covariance matrix $\widehat{G}_{\alpha_1\alpha_2}^{(2)} = G_{\alpha_1\alpha_2}^{(2)} + \Delta\widehat{G}_{\alpha_1\alpha_2}^{(2)}$ can be inverted to second order in fluctuations about the mean covariance matrix:

$$\widehat{G}_{(2)}^{\alpha_1\alpha_2} = G_{(2)}^{\alpha_1\alpha_2} - \sum_{\beta_1,\beta_2\in\mathcal{D}} G_{(2)}^{\alpha_1\beta_1}\Delta\widehat{G}_{\beta_1\beta_2}^{(2)}G_{(2)}^{\beta_2\alpha_2}$$
$$+ \sum_{\beta_1,...,\beta_4\in\mathcal{D}} G_{(2)}^{\alpha_1\beta_1}\Delta\widehat{G}_{\beta_1\beta_2}^{(2)}G_{(2)}^{\beta_2\beta_3}\Delta\widehat{G}_{\beta_3\beta_4}^{(2)}G_{(2)}^{\beta_4\alpha_2}$$
$$+ O\left(\Delta^3\right).$$
(21)

This can then be plugged into the argument for the conditional distribution in Eq. (15), and after Taylor expanding in the fluctuations $\Delta\widehat{G}_{\alpha_1\alpha_2}^{(2)}$, the integral in Eq. (14) evaluates to

$$p\big(z^{(2)}\,|\,\mathcal{D}\big) = \frac{1}{\sqrt{\big|2\pi G^{(2)}\big|^{n_2}}} \exp\left(-\frac{1}{2}\sum_{j=1}^{n_2}\sum_{\alpha_1,\alpha_2\in\mathcal{D}} G^{\alpha_1\alpha_2}_{(2)} z^{(2)}_{j;\alpha_1} z^{(2)}_{j;\alpha_2}\right)\left\{\left[1+O\left(\frac{1}{n_1}\right)\right] + \sum_{i=1}^{n_2}\sum_{\alpha_1,\alpha_2\in\mathcal{D}}\left[O\left(\frac{1}{n_1}\right)\right] z^{(2)}_{i1;\alpha_1} z^{(2)}_{i1;\alpha_2}\right.$$

$$\left. +\frac{1}{8n_1}\sum_{i_1,i_2=1}^{n_2}\sum_{\alpha_1,\dots,\alpha_4\in\mathcal{D}} V^{(\alpha_1\alpha_2)(\alpha_3\alpha_4)}_{(2)} z^{(2)}_{i1;\alpha_1} z^{(2)}_{i1;\alpha_2} z^{(2)}_{i2;\alpha_3} z^{(2)}_{i2;\alpha_4}\right\} + O\left(\frac{1}{n_1^2}\right). \tag{22}$$

A Gaussian is then multiplied by terms organized in powers of $1/n_1$. The contributions to zeroth and first order are contained within curved brackets, and the $O(1/n_1)$ in square brackets within the second term indicates subleading corrections to $G^{(2)}_{\alpha_1\alpha_2}$. By taking the logarithm of the right hand side, the action $S\big(z^{(\ell)}\big)$ is isolated, and constant terms can be absorbed into the partition function, and the action can be expressed as

$$S(z^{(2)}) = \frac{1}{2}\sum_{\alpha_1,\alpha_2\in\mathcal{D}}\left[G^{\alpha_1\alpha_2}_{(2)} + O\left(\frac{1}{n_1}\right)\right]\sum_{i=1}^{n_2} z^{(2)}_{i;\alpha_1} z^{(2)}_{i;\alpha_2}$$

$$-\frac{1}{8}\sum_{\alpha_1,\dots,\alpha_4\in\mathcal{D}}\frac{1}{n_1} V^{(\alpha_1\alpha_2)(\alpha_3\alpha_4)}_{(2)}\sum_{i_1,i_2=1}^{n_2} z^{(2)}_{i1;\alpha_1} z^{(2)}_{i1;\alpha_2} z^{(2)}_{i2;\alpha_3} z^{(2)}_{i2;\alpha_4} + O\left(\frac{1}{n_1^2}\right). \tag{23}$$

The higher-order connected contributions to the moments can also be calculated as terms in the action, and each term will have factors of $1/n_1$ with increasing order. In the large $n_1$ limit, only the lowest-order terms will contribute to the action, and in the infinite limit, the output distribution will become Gaussian, reflecting the behavior predicted by the central limit theorem.

### D.  Recursion relation

The method of derivation from the first hidden layer to the second can be applied to find an output distribution $p\big(z^{(\ell+1)}\,|\,\mathcal{D}\big)$ given the knowledge of $p\big(z^{(\ell)}\,|\,\mathcal{D}\big)$. The output distribution in layer $\ell+1$ can be calculated from $p\big(z^{(\ell)}\,|\,\mathcal{D}\big)$ through the integral

$$p\big(z^{(\ell+1)}\,|\,\mathcal{D}\big) = \int_{-\infty}^{\infty}\left[\prod_{j,\alpha} dz^{(\ell)}_{j;\alpha}\right] p\big(z^{(\ell+1)}\,|\,z^{(\ell)}\big) p\big(z^{(\ell)}\,|\,\mathcal{D}\big). \tag{24}$$

The calculation proceeds identically to the one performed to obtain $p\big(z^{(2)}\,|\,\mathcal{D}\big)$. The conditional distribution $p\big(z^{(\ell+1)}\,|\,z^{(\ell)}\big)$ is defined by

$$p\big(z^{(\ell+1)}\,|\,z^{(\ell)}\big) = \frac{1}{\sqrt{\big|2\pi\widehat{G}^{(\ell+1)}\big|^{n_{\ell+1}}}}$$

$$\times\exp\left(-\frac{1}{2}\sum_{i=1}^{n_{\ell+1}}\sum_{\alpha_1,\alpha_2\in\mathcal{D}} z^{(\ell+1)}_{i;\alpha_1}\widehat{G}^{\alpha_1\alpha_2}_{(\ell+1)} z^{(\ell+1)}_{i;\alpha_2}\right), \tag{25}$$

where the stochastic metric in the $(\ell+1)$-th layer $\widehat{G}^{(\ell+1)}_{\alpha_1\alpha_2}$ is given by

$$\widehat{G}^{(\ell+1)}_{\alpha_1\alpha_2} \equiv C^{(\ell+1)}_b + C^{(\ell+1)}_W\frac{1}{n_\ell}\sum_{j=1}^{n_\ell}\sigma^{(\ell)}_{j;\alpha_1}\sigma^{(\ell)}_{j;\alpha_2}. \tag{26}$$

Again, the metric can be split into a mean and fluctuations, and then the various expansions and the integration can take place to produce distribution $p\big(z^{(\ell+1)}\,|\,\mathcal{D}\big)$ from $p\big(z^{(\ell)}\,|\,\mathcal{D}\big)$, as desired. This generalization of calculating the output distributions also provides a generalization of the action $S\big(z^{(\ell)}\big)$:

$$S\left(z^{(\ell)}\right) \equiv \frac{1}{2}\sum_{i=1}^{n_\ell}\sum_{\alpha_1,\alpha_2\in\mathcal{D}} g^{\alpha_1\alpha_2}_{(\ell)} z^{(\ell)}_{i;\alpha_1} z^{(\ell)}_{i;\alpha_2}$$

$$-\frac{1}{8}\sum_{i_1,i_2=1}^{n_\ell}\sum_{\alpha_1,\dots,\alpha_4\in\mathcal{D}} v^{(\alpha_1\alpha_2)(\alpha_3\alpha_4)}_{(\ell)}$$

$$\times z^{(\ell)}_{i1;\alpha_1} z^{(\ell)}_{i1;\alpha_2} z^{(\ell)}_{i2;\alpha_3} z^{(\ell)}_{i2;\alpha_4} + \dots, \tag{27}$$

where

$$g^{\alpha_1\alpha_2}_{(\ell)} = G^{\alpha_1\alpha_2}_{(\ell)} + O\left(\frac{1}{n_1},\dots\right), \tag{28}$$

$$v^{(\alpha_1\alpha_2)(\alpha_3\alpha_4)}_{(\ell)} = \frac{1}{n_{\ell-1}} V^{(\alpha_1\alpha_2)(\alpha_3\alpha_4)}_{(\ell)} + O\left(\frac{1}{n_1^2},\dots\right). \tag{29}$$

Terms such as $g^{(\ell)}_{\alpha_1\alpha_2}$ and $v^{(\ell)}_{(\alpha_1\alpha_2)(\alpha_3\alpha_4)}$ can be thought of as data-dependent couplings in the action. $G^{(\ell)}_{\alpha_1\alpha_2}$ and $\frac{1}{n_{\ell-1}}V^{(\alpha_1\alpha_2)(\alpha_3\alpha_4)}_{(\ell)}$ are the leading-order contributions to the quadratic and quartic couplings. In the limit of large width, these will be the contributing couplings within the action, and the higher-order moments can be neglected. By setting the hidden layer widths $n_1, n_2, \dots, n_L, n_{\ell_{out}}$

to a common value $n \gg 1$, the leading-order couplings in any layer can be calculated according to the recursion relations

$$G_{\alpha_1\alpha_2}^{(\ell+1)} = C_b^{(\ell+1)} + C_W^{(\ell+1)}\langle\sigma_{\alpha_1}\sigma_{\alpha_2}\rangle_{G^{(\ell)}} + O\left(\frac{1}{n}\right), \quad (30)$$

$$\begin{aligned}
\frac{1}{n}V_{(\alpha_1\alpha_2)(\alpha_3\alpha_4)}^{(\ell+1)} &= \frac{1}{n}\left(C_W^{(\ell+1)}\right)^2[\langle\sigma_{\alpha_1}\sigma_{\alpha_2}\sigma_{\alpha_3}\sigma_{\alpha_4}\rangle_{G^{(\ell)}} \\
&\quad - \langle\sigma_{\alpha_1}\sigma_{\alpha_2}\rangle_{G^{(\ell)}}\langle\sigma_{\alpha_3}\sigma_{\alpha_4}\rangle_{G^{(\ell)}}] \\
&\quad + \frac{1}{n}\frac{\left(C_W^{(\ell+1)}\right)^2}{4}\sum_{\beta_1,\ldots,\beta_4\in\mathcal{D}}V_{(\ell)}^{(\beta_1\beta_2)(\beta_3\beta_4)} \\
&\quad \times \langle\sigma_{\alpha_1}\sigma_{\alpha_2}(z_{\beta_1}z_{\beta_2} - g_{\beta_1\beta_2})\rangle_{G^{(\ell)}} \\
&\quad \times \langle\sigma_{\alpha_3}\sigma_{\alpha_4}(z_{\beta_3}z_{\beta_4} - g_{\beta_3\beta_4})\rangle_{G^{(\ell)}} \\
&\quad + O\left(\frac{1}{n^2}\right). \quad (31)
\end{aligned}$$

In evaluating the Gaussian expectation values in these recursion relations, the couplings in the action will exponentially grow or vanish as the depth $\ell$ increases. To regulate this behavior, a neural network is subjected to a criticality analysis, where values of the initialization hyperparameters $C_W$ and $C_b$ are chosen to facilitate an action where the leading-order quadratic coupling $G_{\alpha_1\alpha_2}^{(\ell)}$ value is fixed with depth in the recursion relation. The values of $C_W$ and $C_b$ may have layer dependence in general, but in satisfying the recursion relations in a network with uniform width $n$, they can be made layer-independent. When this criticality condition is met, the leading-order contributions for higher-order moments gain a factor of the depth $\ell$ to a power equal to that of the order $1/n$ in the expansion. The expansion in $1/n$ then becomes an expansion in $r = \ell_{out}/n$. As an expansion parameter, $r$ governs which connected correlations contribute to the output distribution for the network.

Depending on the application, ANNs perform better as they grow deeper or wider, and so practitioners will increase their network's depth, width or both to suit their problem [30–34]. However, problems like exploding/vanishing gradients and overfitting/underfitting also become more prevalent as these parameters are increased.

From Eqs. (28) and (29), it can be seen that the width $n$ suppresses connected contributions to moments as $n \to \infty$. When the width is infinite, the distributions in each layer become multivariate Gaussian distributions. This is again a manifestation of the central limit theorem (CLT), as the sum over all elements in Eq. (6) combined with the factor of $1/n$ from the distribution of weights results in the preactivations $z_{i,\alpha}^{(\ell)}$ being averaged over all elements in the prior layer. For most neural networks that can be tuned to criticality, the biases will be initially set to zero, and will not spoil the Gaussian behavior in the infinite-width limit where the CLT kicks in.

While increasing width averages out correlations between neurons, the depth $\ell$ introduces new ones through the recursion relations in Eqs. (30) and (31). In analogy to field theory, these moment recursion relations with depth can be thought of as a renormalization group flow for the different "couplings" in the action. Tuning $C_W$ and $C_b$ to criticality describes the values of these initialization hyperparameters that result in the variance being at a fixed point, meaning the coupling is constant with depth.

In order to accomplish the goal of a limited number of degrees of freedom that still describe a realistic ANN, a small ratio of depth to width is desirable. In this limit, the width averages out the higher-order connected contributions from the moments in the distribution, while the depth keeps the lowest-order connected contributions present. This allows the distributions to be "nearly Gaussian", possessing only a few contributing connected moments that can be easily understood and analyzed.

As an expansion parameter, $r$ quantifies the degree of correlation between neurons in a network. This results in three regimes describing the initialized output distribution:

- $r \to 0$, all terms in the action dependent on $r$ vanish, and the output distribution becomes Gaussian (effectively infinite width, CLT kicks in). These networks have turned off their correlations.

- $0 < r \ll 1$, the moments are controllable, truncated, and nontrivial, as is desired in our effective theory approach. These networks are in what is known as the "effectively deep" regime.

- $r \geq 1$, the moments are strongly coupled, and every term in the r expansion contributes to the action. In this regime, the theory becomes highly non-perturbative, and an effective description becomes impossible [8].

With these considerations, ANNFT offers prescriptions for architecture and initialization hyperparameters. To have a network that can be treated as an effective field theory expansion in $r$, the value of $r$ must lie within the effectively deep regime, dictating that the width in the hidden layers must be much larger than the depth. In the case of practitioners increasing width or depth depending on the application, these hyperparameters should be increased together so that $r$ does not enter the non-perturbative or Gaussian regime. Furthermore, the activation function should be chosen such that it admits criticality, with the initialization hyperparameters tuned to critical values. These conditions partially stabilize the network output, and further stabilization comes from training considerations.

## E. Training

All discussion so far has been focused on the initialization of the network, but of course training must also be considered.

The neural tangent kernel (NTK) is the central object for the analysis of network training in ANNFT. The NTK governs the evolution of network outputs during training for ANNs and can depend on the outputs $z_{i,\alpha}$ in not just the final layer, but all layers prior due to the chain rule derivatives present in loss updates during training. The NTK is defined as

$$H^{(\ell)}_{i_1 i_2;\alpha_1 \alpha_2} = \sum_{\mu,\nu} \lambda_{\mu\nu} \frac{dz^{(\ell)}_{i_1;\alpha_1}}{d\theta_\mu} \frac{dz^{(\ell)}_{i_2;\alpha_2}}{d\theta_\nu}, \qquad (32)$$

where $\lambda_{\mu\nu}$ is the learning rate tensor [8]. $\lambda_{\mu\nu}$ is chosen to be a diagonal matrix with separate learning rates for the weights and biases:

$$\lambda_{b^{(\ell)}_{i_1} b^{(\ell)}_{i_2}} = \delta_{i_1 i_2} \lambda^{(\ell)}_b, \quad \lambda_{W^{(\ell)}_{i_1 j_1} W^{(\ell)}_{i_2 j_2}} = \delta_{i_1 i_2} \delta_{j_1 j_2} \frac{\lambda^{(\ell)}_W}{n_{\ell-1}}. \tag{33}$$

This allows for an analysis of the training using separate hyperparameters for weights and biases in a manner similar to $C_W$ and $C_b$ in the analysis of the initialization. Applying this choice to Eq. (32), the NTK becomes

$$\widehat{H}^{(\ell)}_{i_1 i_2;\alpha_1 \alpha_2} = \sum_{\ell'=1}^{\ell} \Bigg[ \sum_{j=1}^{n_{\ell'}} \Bigg( \lambda^{(\ell')}_b \frac{dz^{(\ell')}_{i_1;\alpha_1}}{db^{(\ell')}_j} \frac{dz^{(\ell')}_{i_2;\alpha_2}}{db^{(\ell')}_j} + \frac{\lambda^{(\ell')}_W}{n_{\ell'-1}} \sum_{k=1}^{n_{\ell'-1}} \frac{dz^{(\ell')}_{i_1;\alpha_1}}{dW^{(\ell')}_{jk}} \frac{dz^{(\ell')}_{i_2;\alpha_2}}{dW^{(\ell')}_{jk}} \Bigg) \Bigg]. \quad (34)$$

It should be noted that the NTK is a stochastic quantity for our network, and so a hat is added to $H^{(\ell)}_{i_1 i_2;\alpha_1 \alpha_2}$ to indicate this, like it was for $\widehat{G}^{(\ell)}_{\alpha_1 \alpha_2}$. The terms in square brackets refer to the contribution to the NTK for each layer $\ell' = 1$ to $\ell' = \ell$.

Exploiting the similarity between this quantity and the moments of the pre-training distribution, the NTK's dependence on width and depth is analyzed in the same way. The infinite-width limit for ANNs yielding Gaussian processes results in them being subject to kernel methods. The NTK's training dynamics become governed by a linear differential equation in the infinite-width limit [8, 28, 35–37]. As we back off the infinite-width limit, algorithm-dependent non-linear training dynamics become present, like connected contributions of higher-order moments of the pre-training distribution [8, 35, 38]. As such, the finite-width kernel $H^{(\ell)}$ contains corrections to the linear infinite-width kernel $H^{\{0\}(\ell)}_{\alpha_1 \alpha_2}$ suppressed by powers of $1/n$:

$$H^{(\ell)}_{\alpha_1 \alpha_2} = H^{\{0\}(\ell)}_{\alpha_1 \alpha_2} + \frac{1}{n_{\ell-1}} H^{\{1\}(\ell)}_{\alpha_1 \alpha_2} + \frac{1}{n^2_{\ell-1}} H^{\{2\}(\ell)}_{\alpha_1 \alpha_2} + O\left(\frac{1}{n^3}\right). \qquad (35)$$

The NTK also features depth dependence similar to the couplings of the pre-training distribution. By comparing the expressions for $\widehat{H}^{(\ell)}_{i_1 i_2;\alpha_1 \alpha_2}$ and $\widehat{H}^{(\ell+1)}_{i_1 i_2;\alpha_1 \alpha_2}$ from Eq. (34), the recursion relation for the NTK can be derived:

$$\widehat{H}^{(\ell+1)}_{i_1 i_2;\alpha_1 \alpha_2} = \sum_{j=1}^{n_{\ell+1}} \Bigg( \lambda^{(\ell+1)}_b \frac{dz^{(\ell+1)}_{i_1;\alpha_1}}{db^{(\ell+1)}_j} \frac{dz^{(\ell+1)}_{i_2;\alpha_2}}{db^{(\ell+1)}_j} + \frac{\lambda^{(\ell+1)}_W}{n_\ell} \sum_{k=1}^{n_\ell} \frac{dz^{(\ell+1)}_{i_1;\alpha_1}}{dW^{(\ell+1)}_{jk}} \frac{dz^{(\ell+1)}_{i_2;\alpha_2}}{dW^{(\ell+1)}_{jk}} \Bigg) + \sum_{j_1,j_2=1}^{n_\ell} \frac{dz^{(\ell+1)}_{i_1;\alpha_1}}{dz^{(\ell)}_{j_1;\alpha_1}} \frac{dz^{(\ell+1)}_{i_2;\alpha_2}}{dz^{(\ell)}_{j_2;\alpha_2}} \widehat{H}^{(\ell)}_{j_1 j_2;\alpha_1 \alpha_2}. \quad (36)$$

The NTK analysis then follows a route similar to that of $\widehat{G}^{(\ell)}_{\alpha_1 \alpha_2}$. In the first hidden layer, the NTK $H^{(1)}_{i_1 i_2 \alpha_1 \alpha_2}$ is deterministic like $G^{(1)}_{\alpha_1 \alpha_2}$. This stems from the fact that in the recursion relation Eq. (36), $\widehat{H}^{(0)}_{j_1 j_2;\alpha_1 \alpha_2} = 0$, as the inputs to the network are not being altered by training. Calculating the derivatives in the recursion relation yields

$$H^{(1)}_{i_1 i_2 \alpha_1 \alpha_2} = \delta_{i_1 i_2} \Bigg[ \lambda^{(1)}_b + \lambda^{(1)}_W \Bigg( \frac{1}{n_0} \sum_{j=1}^{n_0} x_{j;\alpha_1} x_{j;\alpha_2} \Bigg) \Bigg], \tag{37}$$

which does not depend on any stochastic variables.

The NTK in all other layers will have stochastic properties that require it to be split into its mean and fluctuations $\widehat{H}^{(\ell)}_{i_1 i_2;\alpha_1 \alpha_2} = H^{(\ell)}_{i_1 i_2;\alpha_1 \alpha_2} + \Delta \widehat{H}^{(\ell)}_{i_1 i_2;\alpha_1 \alpha_2}$. The recursion relations for the different parts of the NTK are calculated in full in chapters 8 and 9 of Ref. [8], but all of the relations depend on architecture hyperparameters such as $C_W$ and $C_b$, and the choice of activation function in addition to the weight and bias learning rates $\lambda^{(\ell)}_W$ and $\lambda^{(\ell)}_b$. Like the couplings of the pre-training distribution, the NTK will exhibit exponential behavior without specific treatment of the hyperparameters. As the architecture hyperparameters are already tuned to criticality, the remaining tuneable hyperparameters are the learning rates.

To avoid this exponential behavior that contributes to the exploding and vanishing gradients, the learning rates are scaled by factors of the width and depth that lead to corrections to the infinite-width NTK $H^{\{0\}(\ell)}_{\alpha_1 \alpha_2}$ controlled by powers of $r$. As an example, the critically scaled learning rates for the ReLU and Tanh activation functions are

$$\lambda^{(\ell)}_{b,\text{ReLU}} = \frac{\tilde{\lambda}_{b,\text{ReLU}}}{\ell_{out}}, \quad \lambda^{(\ell)}_{W,\text{ReLU}} = \frac{\tilde{\lambda}_{W,\text{ReLU}}}{\ell_{out}}, \qquad (38)$$

$$\lambda^{(\ell)}_{b,\text{Tanh}} = \frac{\tilde{\lambda}_{b,\text{Tanh}}}{\ell}, \quad \lambda^{(\ell)}_{W,\text{Tanh}} = \tilde{\lambda}_{W,\text{Tanh}}. \qquad (39)$$

The learning rates for the weights used in practice also include a factor of $1/n$ (and $1/n_0$ in the first layer) to cancel sums over the width in the recursion relations similar to the scaling of $C_W$. These scalings adjust the learning to provide stability to training, and alter the size of the

"raw" rates $\tilde{\lambda}_W$ and $\tilde{\lambda}_b$ to allow for equivalence of learning rates across the different layers and between different network architectures.

This prescription is applied to networks using stochastic gradient descent learning algorithms, which feature a fixed learning rate. In conjunction with critically tuned initialization, this is claimed to result in stable networks that avoid the exploding and vanishing gradient problem. In addition, these networks and their behavior before and after training can be analyzed and understood using techniques from perturbation theory in QFT and the renormalization group, as promised by ANNFT.

In the next two sections, many of these claims will be validated using two-input (the neutron number $N$ and the proton number $Z$) ANNs trained by data from the AME2020 binding energy dataset.

## III. VALIDATING ANNFT INITIALIZATION DISTRIBUTION BEHAVIOR

The first claim we consider is that the large-width limit of the neural networks tends toward a GP prior for the output distribution of functions. We know mathematically that this must happen eventually; the relevant question is how rapidly the limit is approached. We test this first by initializing the network at various widths, with a fixed depth and fixed input. Two different activation functions were used when generating these distributions, ReLU and Softplus. The ReLU activation function is one of the most popular activation functions used in neural networks. ReLU is zero for negative inputs, and linear for inputs greater than or equal to zero. The Softplus activation function is designed to be a smoothed ReLU, but has been observed to be less stable during training [39–41]. These two activation functions were chosen for their similarity and because ReLU possesses a fixed point in its variance, whereas Softplus does not.

Figure 4 shows the histograms of $500,000$ total outputs from different initializations for the ReLU and Softplus activation functions with a fixed depth of two hidden layers and widths from 1 to 1000 neurons per layer. $^{56}$Fe was used as the test input to generate all these distributions, and all inputs to the network were scaled according to a MinMax scaler. In applying ANNFT, inputs to the network should be scaled to be approximately $\mathcal{O}(1)$ [8]. MinMax scaling converts the inputs to be within a range 0 to 1 using the transformation

$$X_{\text{scaled}} = \frac{X - X_{\min}}{X_{\max} - X_{\min}}, \quad (40)$$

on both the number of neutrons $N$ and protons $Z$, with the min and max terms referring to the largest and smallest values of the dataset. The red line in Fig. 4 shows a normal distribution using the mean and variance of the given histogram to compare the output distribution to what the distribution would look like if it were Gaussian. As the width increases (and $r$ decreases), the output

distributions approach Gaussian distributions at about a width of 120 neurons as the large width averages out connected contributions from higher-order moments in the distribution. Both activation functions' output distributions converge to Gaussian distributions in the infinite-width limit, regardless of whether the activation function has a fixed point in its variance or not. Indeed, the Softplus distributions approach the Gaussian limit somewhat faster than the ReLU distributions.

To fully show that neural network output distributions become GPs as the ratio of depth to width becomes small, we must also look at correlations between inputs. Figure 5 compares corner plots for outputs from a pre-training ReLU-based ANN with scalar inputs ranging from $-1.5$ to $+1.5$. These inputs, being $\mathcal{O}(1)$ or less, are left untreated by the MinMax scaler. The width for all plots is fixed at 120 neurons while the depth increases from one hidden layer up to eight hidden layers. The diagonal histograms are consistent with the approach to Gaussian outputs seen in Fig. 4 (e.g., a clear Gaussian shape for one layer with $r = 1/60$ while a clear non-Gaussian shape for eight layers with $r = 1/15$).

The behavior of the off-diagonal joint distributions is likely unfamiliar to practitioners familiar with the most commonly used GPs (such as RBF or Matern kernels). We can understand the pattern by applying the equations presented in Sec. II. Indeed, kernels for the limiting GPs can be calculated by recognizing that the output distributions are mean zero, and that the variance is the only surviving moment [9]. The variance $G^{(\ell)}$ of outputs $z^{(\ell)}$ for a given layer can be used to find the variance of the next using the equation.

$$G^{(\ell+1)} = C_b + C_W \int_{-\infty}^{\infty} dz^{(\ell)} \, \sigma^2(z^{(\ell)}) \, e^{-\frac{(z^{(\ell)})^2}{2G^{(\ell)}}}. \quad (41)$$

Using this expression, the final layer variance $G^{(\ell_{out})}$ can be calculated, and the output distribution for the GP prior completely determined. To acquire the GP kernel for the prior, all that is required is the covariance of the outputs in the final layer $\langle z^{(\ell)}(x) z^{(\ell)}(x') \rangle$, where $x$ and $x'$ are inputs to the network.

The kernel for an infinitely wide neural network with one hidden layer and a critically initialized ReLU activation function can be calculated analytically [9]:

$$\left\langle z^{(\ell=1)}(x) z^{(\ell=1)}(x') \right\rangle = \frac{2}{\pi} |x \cdot x'| (\pi - \phi) \cos(\phi), \quad (42)$$

$$\phi = \arccos\left(\frac{x \cdot x'}{|x \cdot x'|}\right). \quad (43)$$

These equations impose that the correlation coefficients are either 0 or 1, depending on whether $x \cdot x'$ is less than or greater than zero. This behavior, which is that of a manifestly non-stationary GP, is evident in Fig. 5(a). As the ratio of depth to width increases, correlations grow, as seen in the other panels.

Returning to the pre-training binding energy output distributions, it is important to see that the couplings
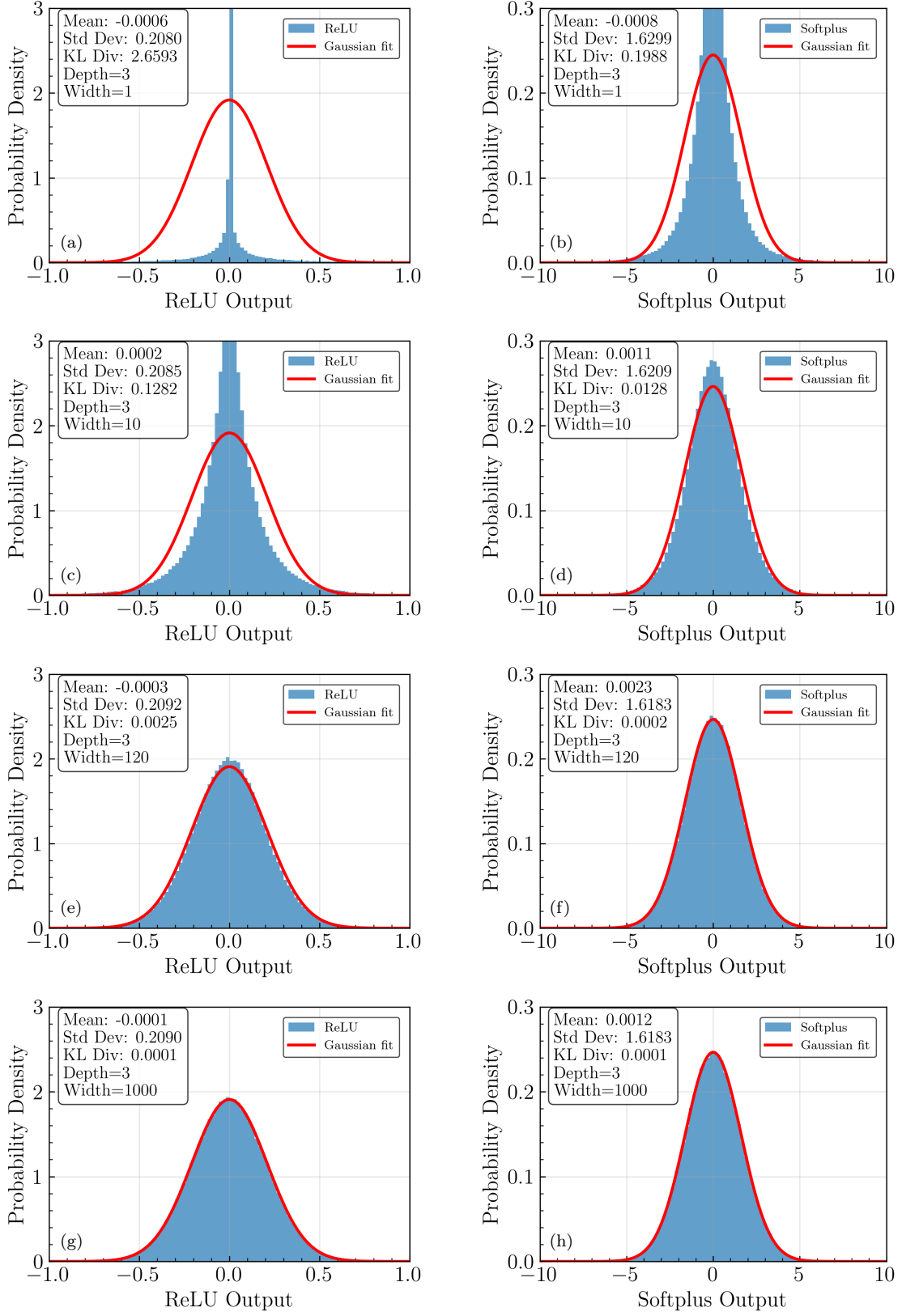
FIG. 4. Final layer ANN output distributions for the scaled $^{56}$Fe input ($Z = 26 \rightarrow 0.13$ and $N = 30 \rightarrow 0.16$) with two hidden layers. Results are given for a ReLU activation function with widths (a) 1, (c) 10, (e) 120, and (g) 1000 neurons per layer and for a Softplus activation function with widths (b) 1, (d) 10, (f) 120, and (h) 1000 neurons per layer.
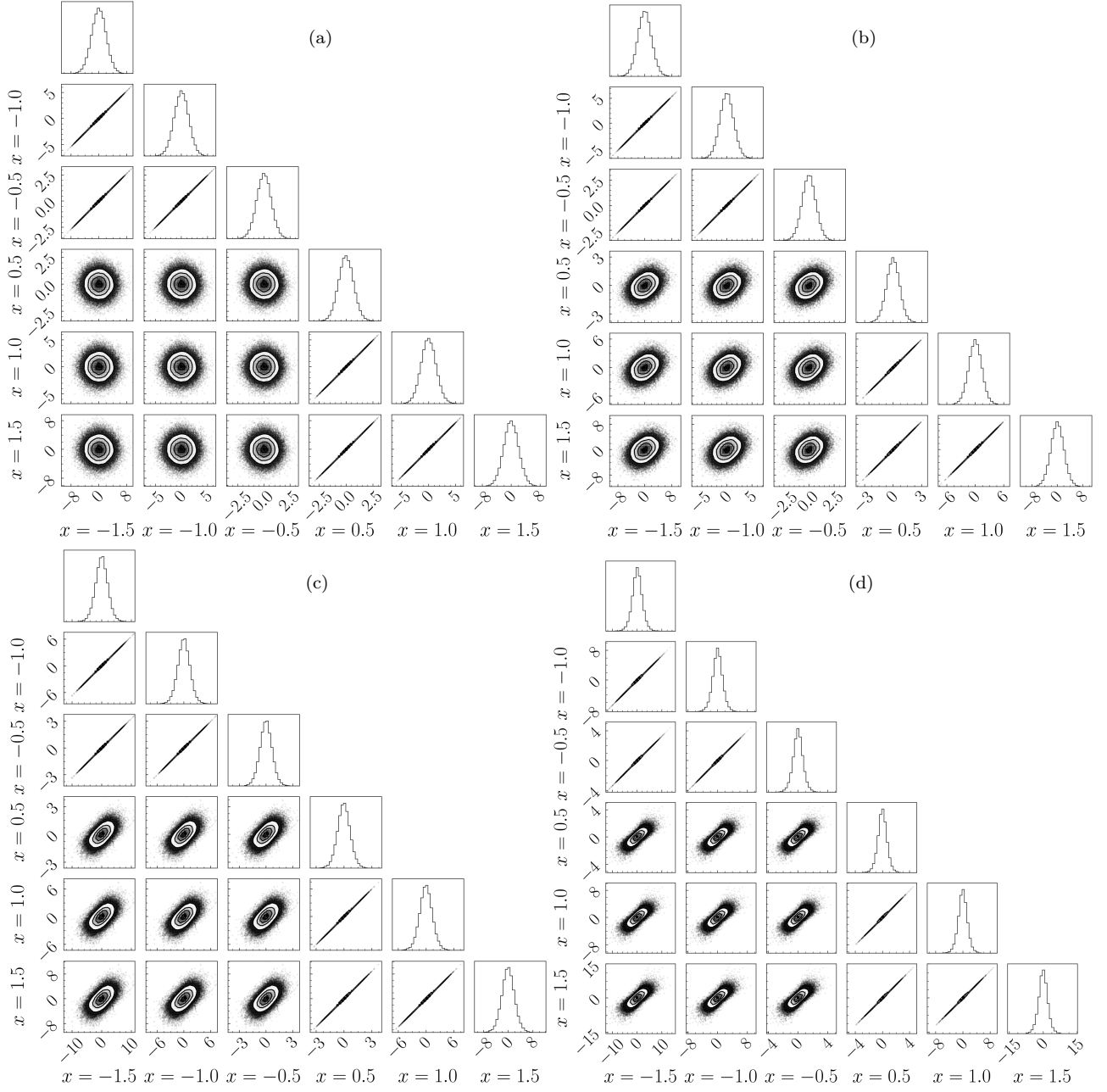
FIG. 5. Corner plots for ReLU activation function outputs. The corner plots all feature ReLU networks of a fixed width of 120 neurons, and depths of (a) 1, (b) 2, (c) 4, and (d) 8. The addition of hidden layers introduces correlations to the network output distributions.
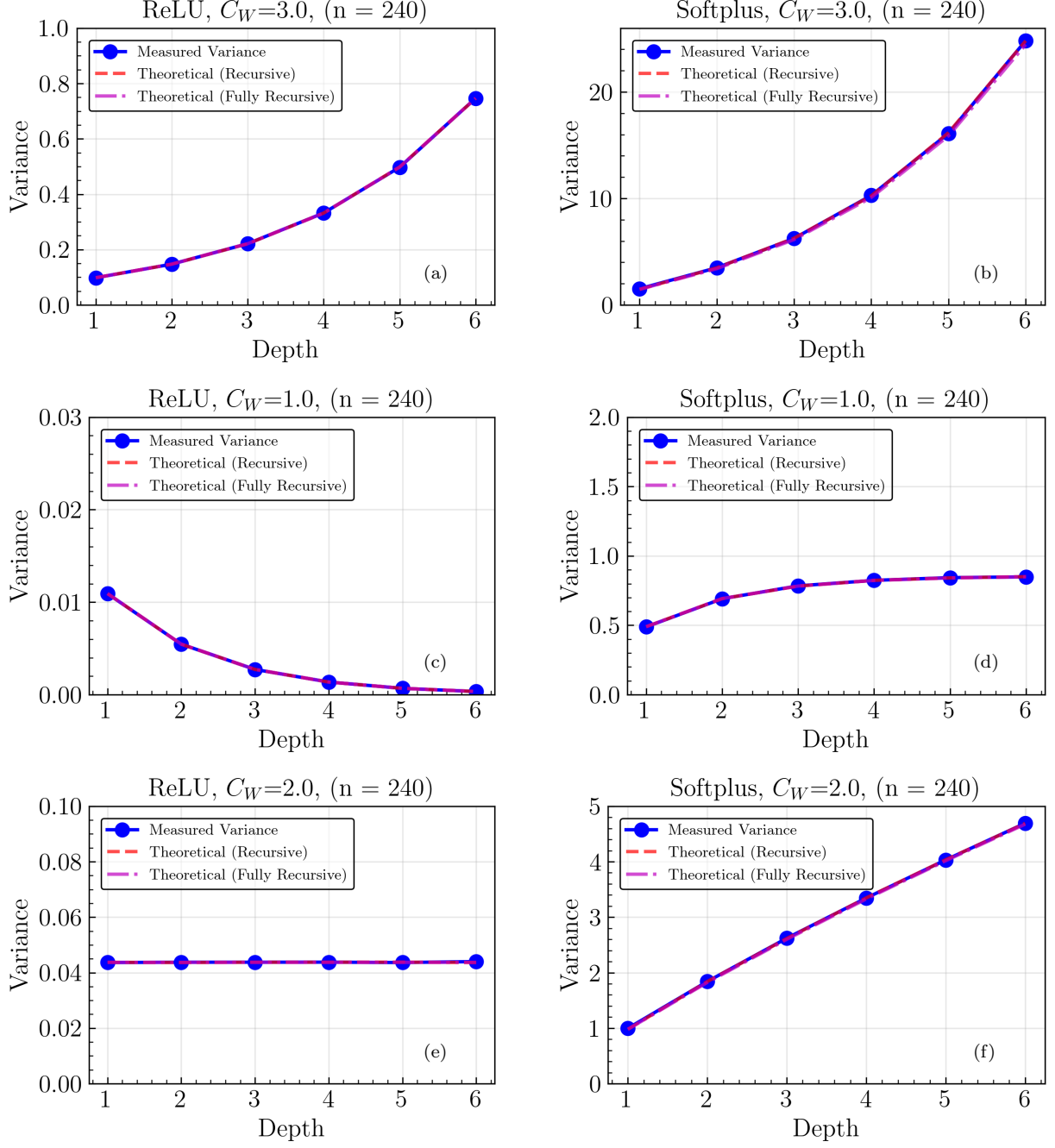
FIG. 6. The final layer pre-training output variance as a function of neural network depth with a fixed width of 240. Results are given for a ReLU activation function tuned (a) above ($C_W = 3.0$, $C_b = 0.0$) (c) below ($C_W = 1.0$, $C_b = 0.0$), and (e) at the critical point ($C_W = 2.0$, $C_b = 0.0$), and for a Softplus activation function with (b) highest, (d) lowest, and (f) intermediate initialization widths for the weights by using the same initialization hyperparameters as the ReLU. The measured variance is plotted in blue, a recursive calculation using empirical values of the variance is plotted as a dashed red line, and a recursive calculation using only theoretically calculated variances is plotted as a dashed pink line. When above/below the critical values of the initialization width, the variance explodes/vanishes with depth. When the ReLU network is tuned to critical initialization, the variance is fixed with depth. The Softplus activations do not have a critical point for initialization. As such, the Softplus variance either grows or asymptotes towards a constant value with depth, and does not have initialization hyperparameters that give a fixed variance with depth.
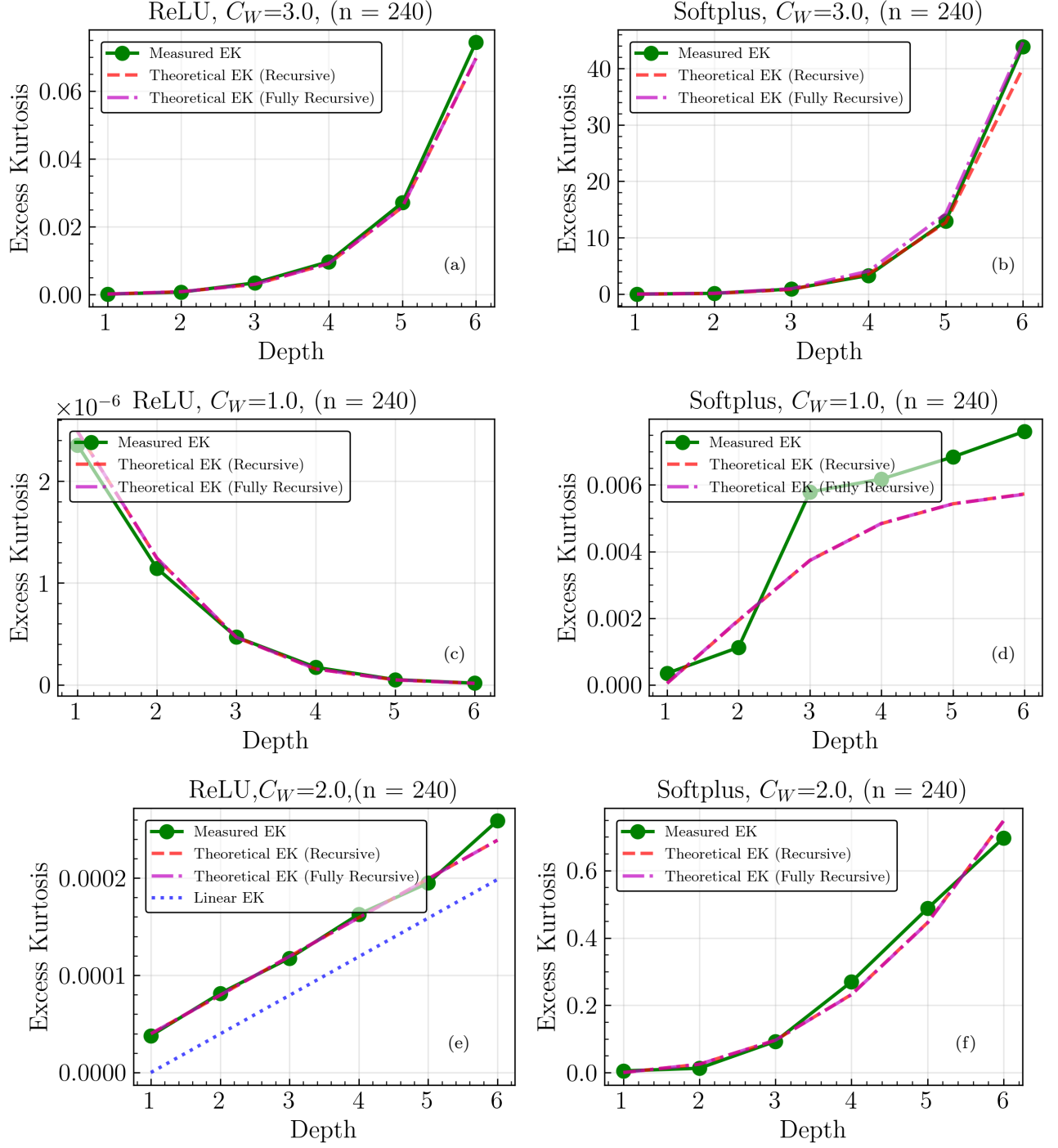
FIG. 7. The final-layer pre-training-output (unstandardized) excess kurtosis (abbreviated here as EK, and also known as the non-Gaussian 4-point correlations) as a function of neural network depth with a fixed width of 240. Results are given for a ReLU activation function tuned (a) above ($C_W = 3.0$, $C_b = 0.0$), (c) below ($C_W = 1.0$, $C_b = 0.0$), and (e) at the critical point ($C_W = 2.0$, $C_b = 0.0$), and for a Softplus activation function with (b) highest, (d) lowest, and (f) intermediate initialization widths for the weights by using the same initialization hyperparameters as ReLU. The measured EK is plotted as a green line, a recursive calculation of EK using empirical values of the variance and initial EK is plotted as a dashed red line, and a recursive calculation using only theoretical values is plotted as a dashed pink line. An additional blue line is present in (e), representing an exact calculation of the critical EK without treating $1/n$ terms as subleading. When critically tuned, ReLU networks' EK behaves linearly with depth as predicted by ANNFT, demonstrating that the higher-order moments of the distribution are controlled by an expansion in $r$. Non-critical distributions are seen to grow non-linearly with depth, reflecting that criticality allows for a perturbative ANNFT that lets deeper network network behavior to be analyzed.

in the distributions' actions evolve according to an RG flow with each layer of the ANN. As we have emphasized, finding a fixed point in the variance prevents gradient issues during training. ReLU is an example of an activation function that has a fixed point. In Fig. 6, the variance of the final layer output distribution either explodes (Fig. 6(a)) or vanishes (Fig. 6(c)) with depth when the initialization widths are away from criticality, and remains at a fixed value when tuned to criticality (Fig. 6(e)).

The variance recursion relation given in Eq. (30) can be used to calculate the variance in any layer $\ell$. The variance value used to calculate the initial value for $G^{(1)}$ is $G^{(0)} = C_W(\frac{1}{n_0} \sum_{i=1}^{n_0} x_i^2) + C_b$. Applying the MinMax scaling to a test input of $^{56}$Fe, two calculations of the variance were performed. The first used the measured value of $G^{(1)}$ from the data, then calculated the variance in all further layers. The second used $G^{(0)}$ to calculate $G^{(1)}$ instead, recursively computing all variances in every layer with only theoretically predicted values of $G^{(\ell)}$.

As can be seen in Figs. 6(a), 6(c), and 6(e), there is agreement between the data and both ANNFT approaches for calculating the variance in all layers. ANNFT accurately calculates the variance for neural networks using multiple different initialization hyperparameter values. In addition, the fixed point that ANNFT predicts for the ReLU activation function is manifested in 6(e) with a value of $G^{(0)} = 4.37 \times 10^{-2}$. The presence of this fixed point is significant, as it demonstrates that there exists an initialization that stabilizes the variance with depth. The other initializations lead to the variance of initial neural network outputs growing or shrinking exponentially, which in turn causes the initial ANN outputs to grow or shrink. Effective ANN training is dependent on the size of the outputs, and the fixed point in the variance helps to stabilize gradients.

By contrast, the Softplus activation function does not have critical initialization and thus no fixed point for its variance. In Fig. 6(b), the Softplus's output variance exhibits exponential growth like the ReLU does when $C_W = 3.0$. However, Softplus's variance grows to much larger values than ReLU, which in turn would cause more unstable gradients with increased network depth. Figures 6(d) and 6(f) show the variance as a function of depth becoming concave compared to the convex behavior in Fig. 6(b). With $C_W = 2.0$, the variance for Softplus is still growing with larger values than the ReLU distribution, but the growth is linear. When $C_W = 1.0$, the variance appears to be approaching a constant value with depth. To arrive at a true fixed point, $C_W$ must approach 0, which is not possible as this choice of initialization would not break the permutation symmetry required to have a trainable neural network [8]. For an activation function to have a fixed point in their variance they need to meet the conditions:

$$\sigma(0) = 0, \tag{44}$$
$$\sigma'(0) \neq 0, \tag{45}$$

and Softplus does not meet the first condition.[4] Despite being intended as a smooth replacement for the ReLU activation function, Softplus is much more prone to exploding/vanishing gradients and therefore inferior to ReLU [39–41].

The next important moment for ANN output distributions is the fourth moment, or kurtosis. The excess kurtosis (EK), which is the difference of the full kurtosis from the product of covariances (so EK is zero for a Gaussian distribution):

$$EK = \langle z_{i_1;\alpha_1} z_{i_2;\alpha_2} z_{i_3;\alpha_3} z_{i_4;\alpha_4} \rangle$$
$$- \langle z_{i_1;\alpha_1} z_{i_2;\alpha_2} \rangle \langle z_{i_3;\alpha_3} z_{i_4;\alpha_4} \rangle, \tag{46}$$

is given by the connected four-point interaction in the ANNFT formulation.[5] In the infinite-width limit, the connected contributions from higher-order moments are suppressed and the neural network becomes a Gaussian process with training determined by a linear differential equation. These higher-order moments' connected contributions allow for more complex correlations between neurons and network's abilities to represent complex features in data. As noted in Sec II, tuning to criticality allows for the action of a neural network to be treated through an expansion in the ratio of depth to width $r$, and the connected contributions to the distribution to be treated as perturbations to the infinite-width Gaussian distribution, or free field theory.

The second term in an action for an ANN with mean zero Gaussian initialization distributions is $V^{(\ell)}/n$, which is the EK of the distribution without the normalization by the variance squared. The quantity $V^{(\ell)}$ satisfies a recursion relation

$$V^{(\ell+1)} = \frac{C_W}{G^{(\ell)}} \langle z\sigma(z)\sigma'(z) \rangle_{G^{(\ell)}} V^{(\ell)}$$
$$+ C_W^2 [\langle \sigma^4(z) \rangle_{G^{(\ell)}} - \langle \sigma^2(z) \rangle_{G^{(\ell)}}^2]. \tag{47}$$

If left untreated, $V^{(\ell)}$ will cause the EK to grow and shrink exponentially with depth. However, with a critically tuned variance $G^{(\ell)}$, the excess kurtosis grows linearly in $\ell$ (and therefore in $r$ as well).

This behavior can be seen for the ReLU activation function in Figs. 7(a), (c), and (e). Along with the empirical EK from the generated output distributions, two calculations of the kurtosis were performed using the recursion relation in Eq. (47). The first computed $V^{(1)}$

---

[4] Note that Softplus *can* meet both conditions for criticality by adding a constant $C = -\log(2)$ so that $\sigma(0) = 0$. Not all activation functions without a fixed point in the variance can be adjusted in this way, but these criteria provide a way to either adjust an activation function so it can be tuned to criticality, or to disregard an activation function if it cannot be adjusted to meet these conditions.

[5] Kurtosis is typically defined as a standardized fourth moment, meaning it is divided by the product of covariances. Our definition is not standardized.

directly from the distribution using Scipy's stats.kurtosis function to find the excess kurtosis, and then scaling the value appropriately to acquire $V^{(1)}$ [42]. From there, the $V^{(\ell)}$s in all subsequent layers were computed recursively using the $G^{(\ell)}$ calculated directly from the output distributions. The second calculation used initial values of $G^{(0)} = C_W(\frac{1}{n_0}\sum_{i=1}^{n_0} x_i^2) + C_b$ and $V^{(0)} = 0$ (the output is Gaussian from Sec II) to compute $V^{(\ell)}$ completely from theoretically predicted values.

Figures 7(a) and (c) show the excess kurtosis growing or shrinking exponentially in agreement with both versions of the recursive calculation. Figure 7(e) shows the desired linear behavior from criticality, and closely matches both recursive calculations. Figure 7(e) also features an additional line corresponding to an exact calculation of $V^{(\ell)}/n$ to leading-order in $1/n$ calculation from Eq. (47) for ReLU:

$$V^{(\ell)} = 5(G^{(\ell)})^2 \frac{\ell - 1}{n} + \mathcal{O}(\frac{1}{n^2}). \qquad (48)$$

If this expression is rearranged such that the expansion is leading-order in $r$, and the $1/n$ contribution present becomes a subleading contribution to the expansion, then

$$V^{(\ell)} = 5(G^{(\ell)})^2 r + \mathcal{O}(r^2), \qquad (49)$$

and the blue line matches the data and recursive calculations. The recursive calculations show agreement with the data, with the exception of Fig. 7(d) (for which the excess kurtosis is very small and is subject to large relative fluctuations).

In the case of the Softplus activation, Fig. 7(b) and (f) demonstrate the excess kurtosis growing exponentially, while 7(d) exhibits concave behavior with increasing depth. The excess kurtosis asymptotically approaching a constant value for $C_W = 1.0$ rather than being linear means these terms in the action are not being adjusted as depth increased. The additional correlations that grow linearly with depth in an output distribution for a critically tuned activation function allow for more expressivity, and are controlled by the width of the network. With a constant EK, the Softplus activation function becomes less expressive and capable of learning features of a dataset for $C_W = 1.0$. Having a variance that remains constant when tuned to criticality, ReLU's excess kurtosis scales linearly with the ratio of depth to width $r$ [8]. This roughly linear behavior can be seen in Fig. 7(e), and away from criticality, the kurtosis explodes. The kurtosis being linear in $r$ demonstrates the presence of an expansion in $r$, and for small $r$ allows for control of statistical moments order by order in $r$.

## IV. VALIDATING ANNFT TRAINING BEHAVIORS

The authors of Ref. [6] trained their binding-energy neural networks with GeLU activation functions and the Adam optimizer [43]. However, the critical prescription for training discussed in Sec. II E is applicable in practice only to the stochastic gradient descent (SGD) optimizer, so we focus on that and return at the end to consider adaptive optimizers. To test the ANNFT training behaviors for the two-input binding-energy ANN from Ref. [6], we trained networks with ReLU and Tanh activation functions using SGD, as the critical analysis described in Secs. II D and II E finds them to be the most stable for neural network training [8]. For each activation function, different combinations of critical and non-critical initialization and training were considered. In particular, we explored critical-initialization critical-training (CICT), critical-initialization non-critical-training (CINT), and non-critical-initialization non-critical-training (NINT). All training used a mean absolute error (MAE) loss function to offer direct comparison to Ref. [6].

In training their two-input binding-energy networks, Ref. [6] used a learning rate for their Adam optimizer of 0.0001, and decay constants 0.9 and 0.999. To make a fair comparison, the SGD-optimized networks were trained in two phases. In the first phase, the unscaled learning rates for the weights $\tilde{\lambda}_W$ and biases $\tilde{\lambda}_b$ are adjusted to be much larger than 0.0001 to allow the network to initially take much larger steps toward a minimum in the parameters. The second phase switches to a much smaller learning rate value after 5000 epochs. The second learning rate's smaller size is intended to fine-tune the network parameters and match the $\mathcal{O}(10^{-4})$ value of the adaptive case. The CICT networks feature the critically scaled learning rates from Eq. (39). As discussed in Sec. II E, critical training scales the learning rates so they depend on the depth and the width to avoid exploding or vanishing gradients. As such, layers may feature different learning rates depending on the choice of activation function, whereas traditionally a global learning rate is used for training. For the SGD trained networks without layer dependence (CINT and NINT), a global learning rate is calculated by averaging over the layer-dependent learning rates of the CICT ReLU or Tanh network with the same depth, width, and activation function to be used in the CINT and NINT case.

Figures 8 and 9 depict the CICT loss in blue, the CINT loss in red, and the NINT loss in green, all as a function of the number of epochs. For ReLU, the fully critical (CICT) network outperforms the half critical (CINT) and non-critical (NINT) network for all depths. Networks with criticality had loss values orders of magnitude lower than the fully non-critical network. In all cases (except depth 1), CICT network outperforms the other two cases. Critical networks (CICT, CINT) improve with depth, while NINT appears to stagnate after two hidden layers.

The controlled analysis of ANNs relies on the applicability of the truncated Taylor series of Eq. (3), which is an expansion about the initialization parameters $\theta$. Therefore, it is important to verify that the difference between the initial and final parameters of the networks is small to
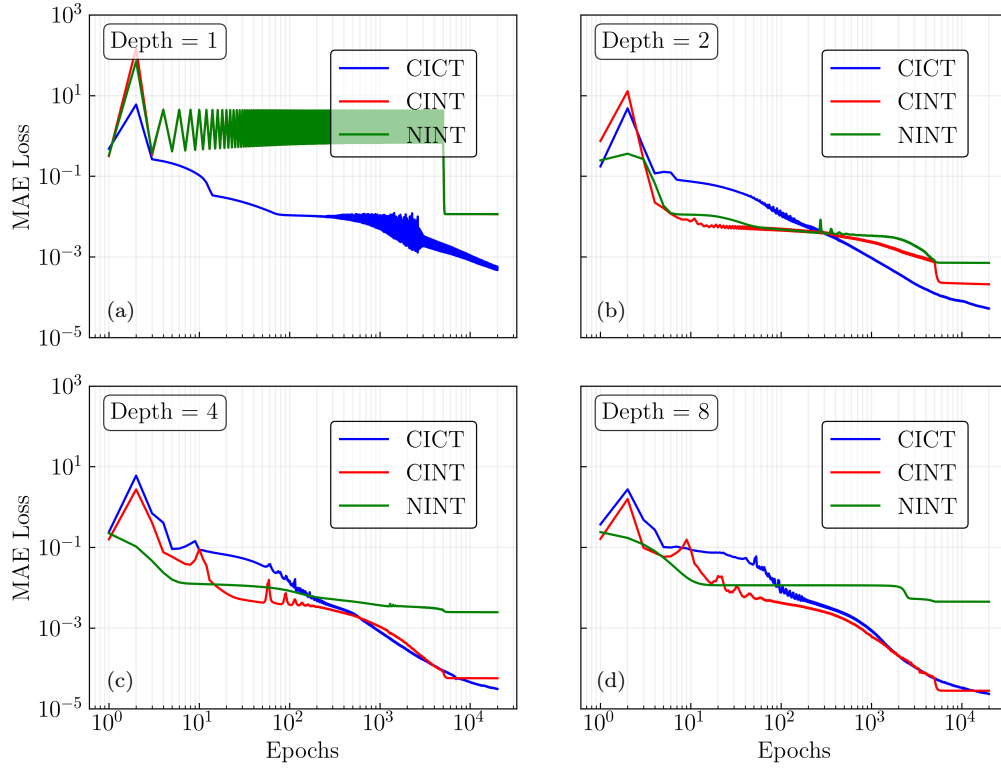
FIG. 8. The mean absolute error loss vs. epochs for ReLU activation functions. Hidden layer widths are at 100 neurons, and depths are 1,2,4, and 8 hidden layers. The CICT and CINT architectures outperform the NINT, and their performance improves with depth, whereas the NINT networks stagnate in performance.
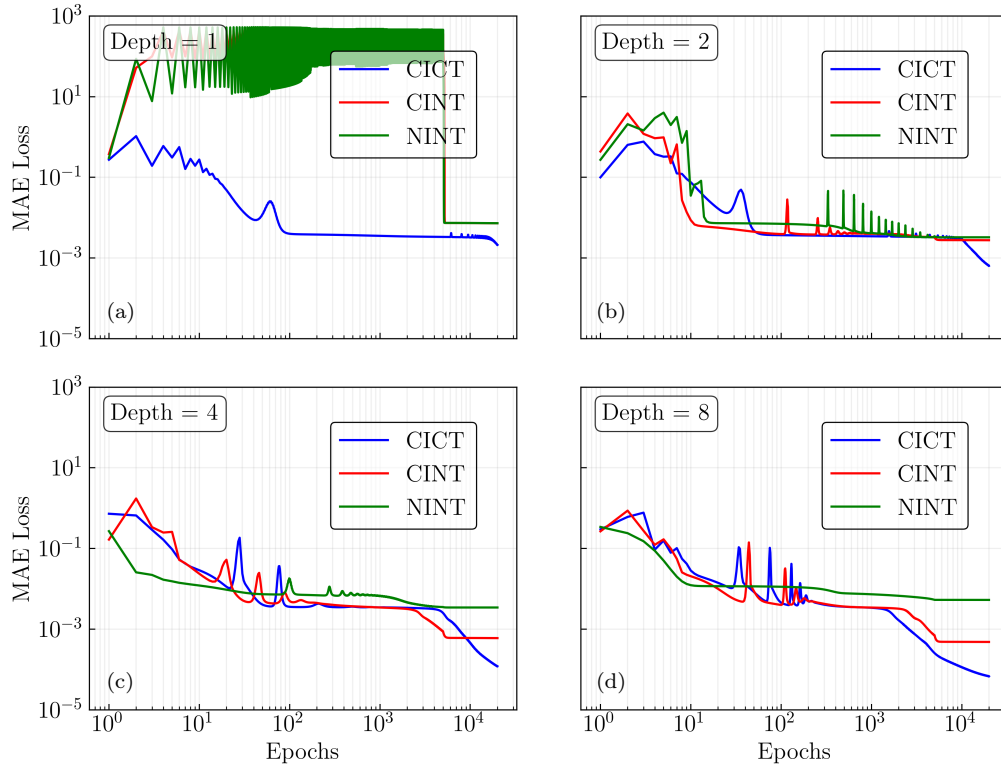


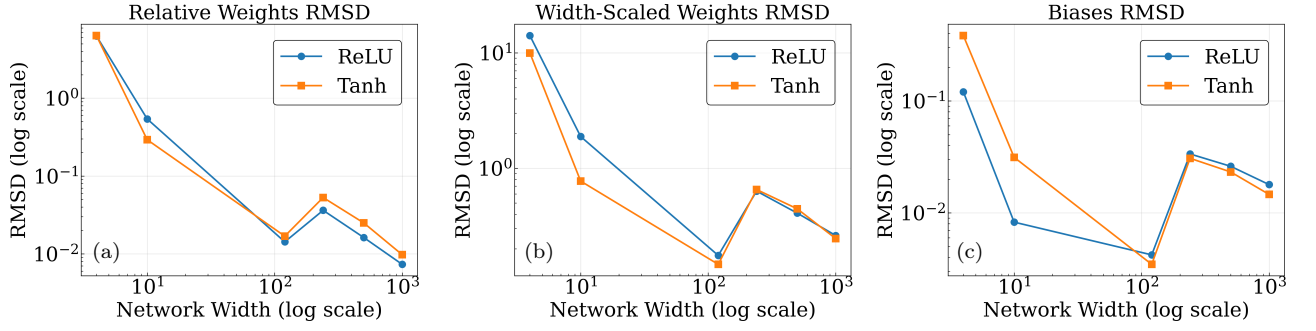FIG. 9. Same as Fig. 8 but for Tanh activation functions.

FIG. 10. Plots of the mean RMSD between final and initial parameters for the average network hidden layer plotted versus multiple widths (4, 10, 120, 240, 500, 1000). (a) depicts the relative RMSD for the weights (compared to the initial weights), while (b) depicts the absolute RMSD with the weight matrices rescaled by a factor of $n$ to remove width dependence, and (c) has the unscaled RMSD for the biases (which are initially zero). Within the regimes useful for training, it can be seen that the changes in parameters are small, a condition required for Taylor expanding around a solution.
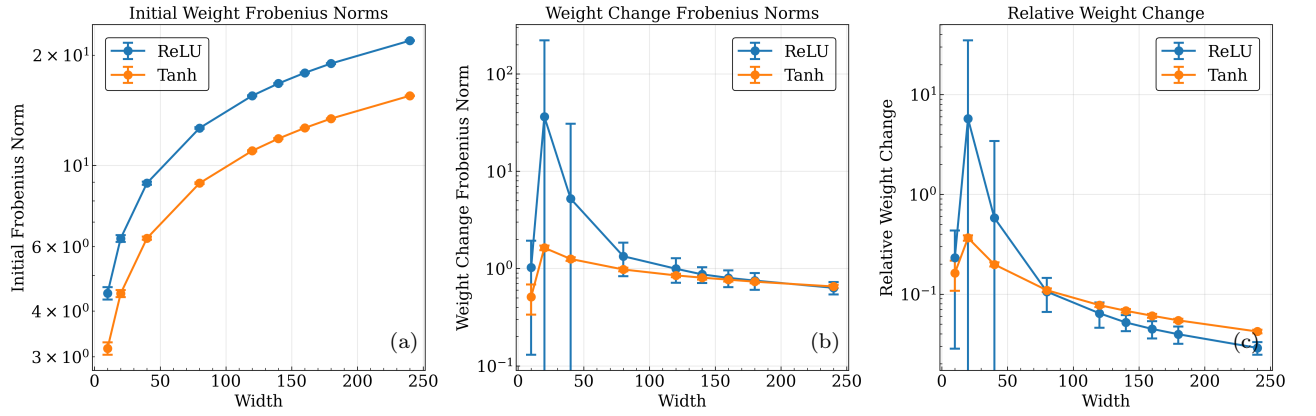


FIG. 11. (a) The mean hidden layer Frobenius norms of the pre-taining weight matrices, (b) the difference between pre and post training matrices, and (c) the relative difference matrices versus the width of the network containing the matrices. Each norm per width is computed by averaging the mean norm for 100 network trainings of 20000 epochs, with outliers outside of the interquartile range of the distribution of each norm being excluded. The depth of each network configuration was held fixed at 4 layers ($L = 3, \ell_{out} = 4$). The norms give an estimate of the size of the initial parameters, and their absolute and relative difference after training, and their dependence on the width of the network.
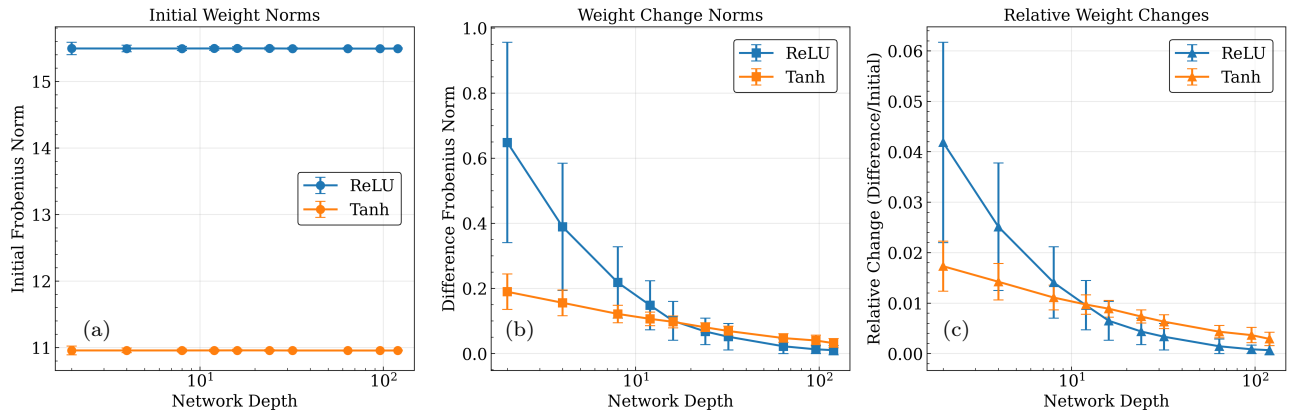


FIG. 12. Same as Fig.11, but now with fixed width instead of depth. (a) The mean hidden layer Frobenius norms of the pre-taining weight matrices, (b) the difference between pre and post training matrices, and (c) the relative difference matrices versus the depth of the network containing the matrices. Each norm per depth is computed by averaging the mean norm for 100 network trainings of 20000 epochs, with outliers outside of the interquartile range of the distribution of each norm being excluded. The width was held fixed at 120 neurons per hidden layer. The norms give an estimate of the size of the initial parameters, and their absolute and relative difference after training, and their dependence on the depth of the network..

ensure that the Taylor series is valid. For critical ReLU and Tanh networks, we find in Figs. 10, 11, and 12 that the difference in the initial and final parameters is indeed small in the regimes of interest for the ratio of depth to width $r$ (see Sec. II D).

Figure 10 depicts the mean of the root-mean-square deviation (RMSD) of the parameter changes in the hidden layers for 100 trainings. The networks used had a fixed depth of $\ell_{out} = 4$ with widths of 4, 10, 120, 240, 500, and 1000. In Fig. 10(a) the mean relative (compared to the initial weights) RMSD of the weights declines rapidly with width $n$. In Fig. 10(b) we show the absolute RMSD with weight matrices scaled by $n$ to normalize the comparison between initial and final weights. We find a similar pattern to the mean relative RMSD in the mean absolute RMSD, which follows because the initialization weights are of order unity. The story for the biases in Fig. 10(c) is similar. This verifies that as the network gets closer to $r \ll 1$, the change in parameters within the network decreases.

This conclusion is seen again in the mean hidden layer Frobenius norms versus widths in Fig. 11. As the width increases, so does the Frobenius norm of the hidden layer matrices in Fig. 11(a), as the matrix is $n \times n$, where $n$ is the width of the network's hidden layers. The difference between initial and final parameters increases with width initially in Fig. 11(b), but decreases for width values greater than 48 neurons. The relative norm in Fig. 11(c) is very small in the region of interest, and its decrease with width shows the desired small parameter changes for fixed-depth networks.

Figure 12 shows the Frobenius norms for networks with a fixed width of $n = 120$ and increasing depth. At fixed width the Frobenius norms remains constant as depth is increased, as seen in Fig. 12(a), as the increased depth will only include new matrices to be averaged over rather than increasing the size of the weight matrices like the fixed-depth case. The absolute changes in difference norms are shown in Fig. 12(b) and the relative changes in Fig. 12(c). At any depth with this large width (up to $r = 1$), the relative changes are small, much smaller than one. In summary, we find that networks with small $r$ values warrant the use of Taylor expansions of the network solutions around their initial parameters $\theta$.

Different ANN learning regimes associated with the value $r$ were identified in Sec. II D (see [8]). When $r \geq 1$, correlations between neurons are too strong to learn features, and the networks can be considered "overly deep". When $r \ll 1$, the networks are "effectively deep", and can learn. When $r \to 0$, the neuron correlations turn off, and features cannot be learned. This implies an optimal $r$ (which will be denoted $r^*$ below). The behavior with $r$ is illustrated in Figs. 13 and 14.

The binding-energy RMSD decreases with $r$ in Fig. 13 until it reaches a minimum at $r = 0.033$. The distribution of RMSD values is roughly a normal distribution up to this point, with widths indicated by the error bars. The decreasing values of the binding-energy
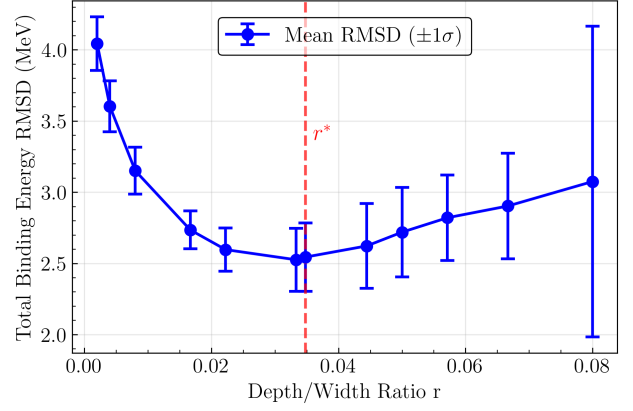


FIG. 13. The mean root-mean-square (rms) deviation of the binding energy data and the two-input-network outputs for 100 trainings vs. ratio of depth-to-width $r$, with the depth $\ell_{out} = 4$. Trainings with binding-energy RMSD$\geq$ 30 MeV were omitted from calculation of the mean value of the binding-energy RMSD. The red horizontal line is the $r^*$ value for this network, calculated from Eq. (50) to be $r^* = 0.034$. This labels the cutoff where the effectively deep regime begins to transition into the chaotic regime with increasing $r$.

RMSD demonstrate the improvement in feature learning as the network moves from the Gaussian limit ($r \to 0$) towards the effectively deep region. Past the minimum value, the distribution of binding-energy RMSD values becomes bimodal, with a second mode centered around about 72 MeV (here the error bars indicate the standard deviation, but a normal distribution is not implied). This other mode is associated with networks that fail to train, corresponding to the transition from effectively deep to overly deep (chaotic), where the network becomes incapable of learning any features. For $r \geq 1$, the distributions become entirely centered around 72 MeV.

The role of $r$ in determining learning capability is visualized in Fig. 14, which shows binding-energy residuals across the nuclear landscape as heat maps for a wide range of $r$ values. In Fig. 14(a) with $r = 0.004$ the residuals are generally small, but an increase to $r = 0.033$ in Fig. 14(b) shows improved learning. A further increase to $r = 0.4$ in Fig. 14(c) shows significant degradation and $r = 1$ in Fig. 14(c) shows poor learning.

Given the minimum binding-energy RMSD seen in Fig. 13, the existence of an optimal value $r$ for a neural network is implied. Roberts and Yaida discuss methods utilizing information theory for calculating $r^*$, the optimal aspect ratio, in appendix A of Ref. [8]. $r^*$ defines the boundary where the effectively deep regime ends, and the chaotic regime begins, and they calculate it (to $\mathcal{O}(r^3)$) to be

$$r^* = \left( \frac{4}{20 + 3n_{\ell_{out}}} \right) \frac{1}{\nu}. \tag{50}$$

Here $n_{\ell_{out}}$ is the size of the output layer ($n_{\ell_{out}} = 1$ for the binding-energy network) and $\nu$ is equal to the slope
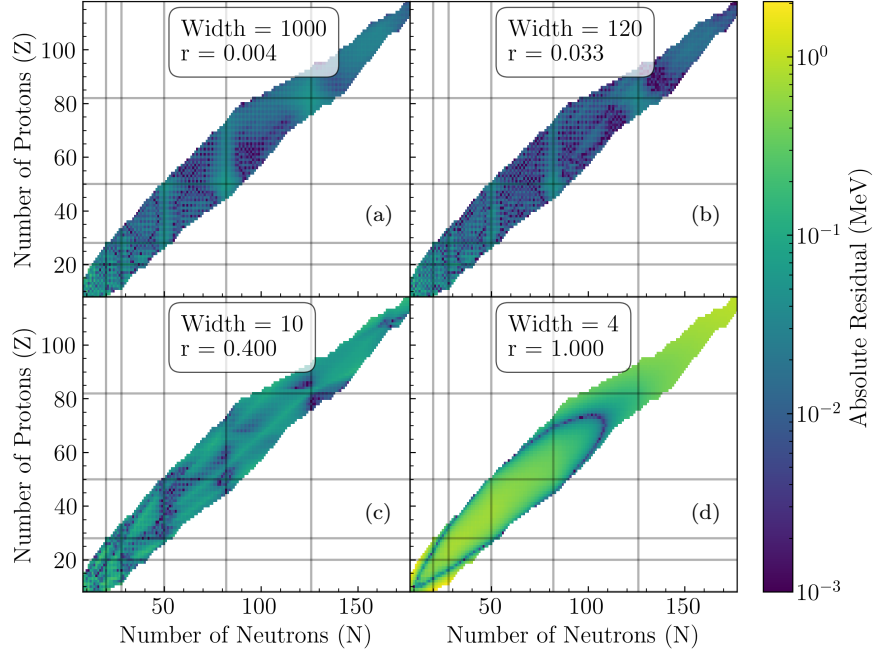
FIG. 14. Residual plots for a trained 2 input binding energy network with a fixed depth of 4, and widths of (a) 4, (b) 10, (c) 120, and (d) 1000 to demonstrate the different learning regimes in ANNFT.

factor of the critical linear EK; $\nu = 5$ for the ReLU network EK at criticality and $r^* = 0.034$. For our $\ell_{out} = 4$ case, this suggests the optimal width of $n = 115$. This is consistent with Fig. 13, which has a minimum binding-energy RMSD when $r = 0.033$, and growing values of the binding-energy RMSD past the $r^*$ cutoff. As noted above, the distributions of outputs for runs with $r > r^*$ feature more non-converged runs. This suggests that while architectures with $r^* < r < 1$ can result in successfully trained networks, they have a greater chance of failing as $r$ increases. Table I displays the width $n$, ratio of depth to width $r$, the mean and std of the binding-energy RMSD, and the number of runs used in calculating the statistics for the RMSD. Runs where the RMSD was greater than 100 MeV, or produced a NaN value are excluded from calculation of the statistics, and 100 runs were used at most for each $r$ value. The mean and std have an asterisk for $r$ values where $r^* < r < 1$, where the distribution of RMSD values are split between the effectively deep and chaotic regimes as noted before. Indeed, one can see the mean and std values transition between $\mathcal{O}(1)$ means and $\mathcal{O}(10^{-1})$ std to much larger means and variances before settling entirely at around 72 MeV when $r = 4$.

To summarize, we find activation functions should be chosen such that critical points exist. Initialization widths for weights and biases should be set to criticality to avoid gradient problems during training. Learning rates must be critically scaled to improve performance and offer further gradient stability. The ratio of depth to width $r$ should be small ($0 < r < r^* \ll 1$) to op-

timally learn the features of the system. Small $r$ also means ANNs can be treated as field theories with running couplings that run with depth, and can be tuned to criticality. $r^*$ provides a more specific constraint for the boundary between the overly deep $r \geq 1$ regime, and the effectively deep regime beyond the condition that $r \ll 1$ [8].

We now return to the use of adaptive optimization as opposed to SGD. To make a comparison, the SGD optimization in the ReLU and Tanh networks is replaced by an Adam optimizer using the same training hyperparameters as Ref. [6], with the exception of the learning rate. The learning rate was calculated instead using the same averaging over the layer-dependent critical learning rates used for the non-critical training systems (CINT and NINT). At present, there is no critical training prescription for the Adam optimizer, and so the only cases tested and compared with adaptive Adam optimizers were CINT and NINT.

The bottom line of the comparison is that when the SGD optimizer was replaced by Adam, both CINT and NINT networks no longer show criticality, but obtain similar binding-energy RMSD values to the original Zeng network of about 1.2 MeV [6]. Thus the adaptive algorithm apparently compensates for non-critical initialization and training. The (non-adaptive) CICT network gets a best binding-energy RMSD of about 1.9 MeV, which is only a small factor worse than the result from the Adam optimizer. Nevertheless, the improved adaptive optimizers offer improved training performance. However, this is at the cost of obscuring the analysis of the

TABLE I. RMSD mean and standard deviation (std) $\sigma$, and the number of runs used for the distributions, for a large range of network widths and ratio of depth to width ($r$). For this table, the RMSD cutoff whenever averaging over runs was increased to 100 MeV. Values of the mean and std where $r^* < r < 1$ are marked with an asterisk; as $r$ increases above $r^*$ the distribution of RMSDs becomes increasingly bimodal.

| Width | Ratio ($r = \ell_{out}/n$) | Mean RMSD (MeV) | $\sigma$ RMSD (MeV) | Number of Runs |
|---|---|---|---|---|
| 1 | 4.0000 | 71.9549 | 0.0000 | 100 |
| 4 | 1.0000 | 68.9716 | 11.2120 | 100 |
| 10 | 0.4000 | 55.9824* | 25.3620* | 99 |
| 20 | 0.2000 | 35.2627* | 31.6383* | 99 |
| 30 | 0.1333 | 14.9702* | 24.1069* | 98 |
| 40 | 0.1000 | 6.0500* | 12.0067* | 99 |
| 50 | 0.0800 | 3.0745* | 1.0921* | 100 |
| 60 | 0.0667 | 2.9027* | 0.3699* | 100 |
| 70 | 0.0571 | 2.8203* | 0.3005* | 100 |
| 80 | 0.0500 | 2.7189* | 0.3137* | 100 |
| 90 | 0.0444 | 2.6223* | 0.2980* | 100 |
| 115 | 0.0348 | 2.5392 | 0.2326 | 100 |
| 120 | 0.0333 | 2.5250 | 0.2221 | 100 |
| 180 | 0.0222 | 2.5964 | 0.1523 | 100 |
| 240 | 0.0167 | 2.7347 | 0.1328 | 100 |
| 500 | 0.0080 | 3.1506 | 0.1641 | 100 |
| 1000 | 0.0040 | 3.6029 | 0.1780 | 100 |
| 2000 | 0.0020 | 4.0420 | 0.1850 | 84 |

neural network. A critical training prescription consistent with adaptive algorithms would be most welcome.

Additional trials were performed for a four-input network, which included the inputs $Z_0$ and $N_0$ to account for pairing effects, following Ref. [6]. $Z_0(N_0)$ is equal to 0 when $Z(N)$ is even and 1 when $Z(N)$ is odd, and as such these inputs were unscaled by the MinMax scaler. The four-input network exhibited the same initialization and training behavior presented in Secs. III and IV. In particular, we find that an increase in input size still possesses the essential properties of ANNFT, and allows for additional features to be added to networks if desired.

Unlike the two-input case, however, the total-dataset binding-energy RMSD values of the four-input network using critical prescriptions were not able to match the RMSD value of Ref. [6]'s optimized four-input network. This will require further investigation. We note that even with a four-input network using the exact architecture and training hyperparameters in Ref. [6] and trained with an adaptive optimizer (Adam) on the same dataset, we were unable to reproduce the same total-dataset binding-energy RMSD values found.

## V. CONCLUSION AND OUTLOOK

Neural network criticality offers a means to understand neural network behavior, as well as providing prescriptions for architecture and training hyperparameters. In this work, we have explored the non-empirical approach to understanding ANNs using a field-theory-based criticality analysis (ANNFT) [8]. We used a protoypical nuclear physics ANN as a testbed, namely a two-input/one-output feed-forward network trained from the AME2020

dataset to fit nuclear binding energies [6].

We first confronted ANNFT predictions of initialization distribution behavior in Sec. III. This is pre-training; we simply collect statistics from repeated random initializations of the weights and biases. We found the expected approach to Gaussian distributions for the final layer output distributions with fixed depth but increasing width. Within small fluctuations, normal distributions were seen by a width of 120 neurons for both ReLU and Softplus activations (Fig. 4) and for every other activation function we have tried. The correlations between pairs of input values with a width of 120 and increasing depths showed Gaussian process (GP) correlations with the predicted kernel at depth 1 and increasing correlations with additional hidden layers (Fig. 5).

A comparison of measured variance and excess kurtosis (difference of the fourth moment from Gaussian expectations) from many initializations for a fixed width of 240 and increasing depths validates the theoretical predictions of criticality (Figs. 6 and 7). In particular, for the ReLU activation function the variance grows or shrinks with depth if the initialization hyperparameters are tuned above or below the predicted critical value. At criticality the variance is flat. In contrast, the Softplus activation function does not have a critical fixed point for initialization. In accordance with theory, the variance either grows or asymptotes to a constant value with depth. Non-zero excess kurtosis (EK) indicates non-Gaussian four-point correlations. For both activation functions, the observed behavior of the EK with depth is in good accord with theoretical expectations.

Next we considered training of the ANN in Sec. IV. We followed the mean absolute error (MAE) loss for a large number of epochs with several different depths, us-

ing combinations of critical and non-critical initialization and training. The networks with critical initialization and critical training (CICT) reached the lowest loss in all cases both for ReLU and Tanh activation functions, with particular contrast to when non-critical hyperparameters were used for both initialization and training (Figs. 8 and 9).

Most practitioners do not encounter these issues even though they rarely pay attention to details of initialization and training. This is because the critically tuned networks are still outperformed by adaptive optimizers. Nevertheless, the results are comparable, and differ by a much smaller factor than when untuned.

Of particular interest is the prediction by ANNFT of different training regimes, characterized by the ratio of depth to width $r$. The predictive analyses of ANNFT are based on the validity of a Taylor series expansion about the initialization. By considering the mean hidden-layer rms deviation and Frobenius norms of the weight matrices before and after training, we can see that the conditions of small deviations from initialization are achieved with sufficiently small $r$ (Figs. 10, 11, and 12). The different learning regimes are manifested by looking at the rms deviation of the binding energy data and ANN outputs as a function of $r$ (Fig. 13), achieved by varying the width at fixed depth, and through heatmaps of global residuals across all $N$ and $Z$ for four different $r$ values (Fig. 14).

In future work, we will continue to investigate critical-

ity in binding-energy fitting with more features/inputs and work toward extending criticality to other training algorithms/hyperparameters. At the same time, we will explore more nuclear systems with critical networks to further understand learning behavior in different problems. The field-theory-based framework offers the possibility of applying approximations and techniques from many-body field theory to extend ANNFT beyond simple perturbation theory. Finally, we will consider an alternative formulation of ANNFT that organizes an expansion with non-Gaussianity not around $1/n$ corrections to the infinite-width limit but through parametric violations of the necessary conditions for the CLT (e.g., independence).

[1] A. Boehnlein, M. Diefenthaler, N. Sato, M. Schram, V. Ziegler, C. Fanelli, M. Hjorth-Jensen, T. Horn, M. P. Kuchera, D. Lee, W. Nazarewicz, P. Ostroumov, K. Orginos, A. Poon, X.-N. Wang, A. Scheinker, M. S. Smith, and L.-G. Pang, Rev. Mod. Phys. 94, 031003 (2022), arXiv:2112.02309 [nucl-th].

[2] G. Carleo, I. Cirac, K. Cranmer, L. Daudet, M. Schuld, N. Tishby, L. Vogt-Maranto, and L. Zdeborová, Reviews of Modern Physics 91 (2019), 10.1103/revmodphys.91.045002.

[3] A. M. Deiana, N. Tran, J. Agar, M. Blott, G. Di Guglielmo, J. Duarte, P. Harris, S. Hauck, M. Liu, M. S. Neubauer, J. Ngadiuba, S. Ogrenci-Memik, M. Pierini, T. Aarrestad, S. Bähr, J. Becker, A.-S. Berthold, R. J. Bonventre, T. E. Müller Bravo, M. Diefenthaler, Z. Dong, N. Fritzsche, A. Gholami, E. Govorkova, D. Guo, K. J. Hazelwood, C. Herwig, B. Khan, S. Kim, T. Klijnsma, Y. Liu, K. H. Lo, T. Nguyen, G. Pezzullo, S. Rasoulinezhad, R. A. Rivera, K. Scholberg, J. Selig, S. Sen, D. Strukov, W. Tang, S. Thais, K. L. Unger, R. Vilalta, B. von Krosigk, S. Wang, and T. K. Warburton, Frontiers in Big Data 5 (2022), 10.3389/fdata.2022.787421.

[4] T. Wolfgruber, M. Knöll, and R. Roth, Phys. Rev. C 110, 014327 (2024), arXiv:2310.05256 [nucl-th].

[5] A. Gnech, B. Fore, A. J. Tropiano, and A. Lovato, Phys.

Rev. Lett. 133, 142501 (2024), arXiv:2308.16266 [nucl-th].

[6] L.-X. Zeng, Y.-Y. Yin, X.-X. Dong, and L.-S. Geng, Phys. Rev. C 109, 034318 (2024), arXiv:2210.02906 [nucl-th].

[7] D. A. Roberts, (2021), arXiv:2104.00008 [hep-th].

[8] D. A. Roberts, S. Yaida, and B. Hanin, The Principles of Deep Learning Theory (Cambridge University Press, 2022) arXiv:2106.10165 [cs.LG].

[9] J. Halverson, A. Maiti, and K. Stoner, Mach. Learn. Sci. Tech. 2, 035002 (2021), arXiv:2008.08601 [cs.LG].

[10] J. Halverson, (2021), arXiv:2112.04527 [hep-th].

[11] S. S. Schoenholz, J. Pennington, and J. Sohl-Dickstein, "A correspondence between random neural networks and statistical field theory," (2017), arXiv:1710.06570 [stat.ML].

[12] O. Cohen, O. Malka, and Z. Ringel, Phys. Rev. Res. 3, 023034 (2021).

[13] D. Bachtis, G. Aarts, and B. Lucini, Phys. Rev. D 103, 074510 (2021), arXiv:2102.09449 [hep-lat].

[14] J. Erdmenger, K. T. Grosvenor, and R. Jefferson, SciPost Phys. 12, 041 (2022), arXiv:2107.06898 [hep-th].

[15] H. Erbin, V. Lahoche, and D. O. Samary, Mach. Learn. Sci. Tech. 3, 015027 (2022), arXiv:2108.01403 [hep-th].

[16] H. Erbin, V. Lahoche, and D. O. Samary (2022) arXiv:2212.11811 [hep-th].

[17] C. Bishop, *Pattern recognition and machine learning*, Vol. 4 (Springer New York, 2006).

[18] Y. S. Abu-Mostafa, M. Magdon-Ismail, and H.-T. Lin, *Learning From Data* (AMLBook, 2012).

[19] K. P. Murphy, *Machine Learning: A Probabilistic Perspective* (The MIT Press, 2012).

[20] P. Mehta, M. Bukov, C.-H. Wang, A. G. Day, C. Richardson, C. K. Fisher, and D. J. Schwab, Physics reports **810**, 1 (2019).

[21] C. M. Bishop and H. Bishop, *Deep learning: Foundations and concepts* (Springer, 2024).

[22] R. Abdulkadirov, P. Lyakhov, and N. Nagornov, Mathematics **11** (2023), 10.3390/math11112466.

[23] L. Liu, H. Jiang, P. He, W. Chen, X. Liu, J. Gao, and J. Han, "On the variance of the adaptive learning rate and beyond," (2021), arXiv:1908.03265 [cs.LG].

[24] R. M. Neal, "Priors for infinite networks," in *Bayesian Learning for Neural Networks* (Springer New York, New York, NY, 1996) pp. 29–53.

[25] C. K. I. Williams, in *NIPS* (1996).

[26] J. Lee, Y. Bahri, R. Novak, S. S. Schoenholz, J. Pennington, and J. Sohl-Dickstein, "Deep neural networks as gaussian processes," (2018), arXiv:1711.00165 [stat.ML].

[27] A. G. de G. Matthews, M. Rowland, J. Hron, R. E. Turner, and Z. Ghahramani, "Gaussian process behaviour in wide deep neural networks," (2018), arXiv:1804.11271 [stat.ML].

[28] A. Jacot, F. Gabriel, and C. Hongler, Proceedings of the 53rd Annual ACM SIGACT Symposium on Theory of Computing (2018).

[29] B. Hanin and M. Nica, "Finite depth and width corrections to the neural tangent kernel," (2019), arXiv:1909.05989 [cs.LG].

[30] F.-L. Fan, R. Lai, and G. Wang, "Quasi-equivalence of width and depth of neural networks," (2022).

[31] A. Álvarez López, A. Slimane, and E. Zuazua, Neural Networks **180** (2024), 10.2139/ssrn.4753244.

[32] A. Radhakrishnan, M. Belkin, and C. Uhler, Proceedings of the National Academy of Sciences **120** (2023), 10.1073/pnas.2208779120.

[33] G. Vardi, G. Yehudai, and O. Shamir, "Width is less important than depth in relu neural networks," (2022), arXiv:2202.03841.

[34] T. Nguyen, M. Raghu, and S. Kornblith, "Do wide and deep networks learn the same things? uncovering how neural network representations vary with width and depth," (2020), arXiv:2010.15327.

[35] A. Jacot, F. Gabriel, and C. Hongler, "Neural tangent kernel: Convergence and generalization in neural networks," (2020).

[36] S. Prince, "The neural tangent kernel," (2024).

[37] L. Weng, "Some math behind neural tangent kernel," (2022).

[38] R. Novak, J. Sohl-Dickstein, and S. S. Schoenholz, "Fast finite width neural tangent kernel," (2022).

[39] X. Glorot, A. Bordes, and Y. Bengio, in *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, Proceedings of Machine Learning Research, Vol. 15, edited by G. Gordon, D. Dunson, and M. Dudík (PMLR, Fort Lauderdale, FL, USA, 2011) pp. 315–323.

[40] A. K. Bhoi, P. K. Mallick, C.-M. Liu, and V. E. Balas, eds., *Bio-inspired Neurocomputing*, Vol. 903 (Springer Singapore, 2021).

[41] S. R. Dubey, S. K. Singh, and B. B. Chaudhuri, "Activation functions in deep learning: A comprehensive survey and benchmark," (2022), arXiv:2109.14545 [cs.LG].

[42] P. Virtanen, R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E. Burovski, P. Peterson, W. Weckesser, J. Bright, S. J. van der Walt, M. Brett, J. Wilson, K. J. Millman, N. Mayorov, A. R. J. Nelson, E. Jones, R. Kern, E. Larson, C. J. Carey, İ. Polat, Y. Feng, E. W. Moore, J. VanderPlas, D. Laxalde, J. Perktold, R. Cimrman, I. Henriksen, E. A. Quintero, C. R. Harris, A. M. Archibald, A. H. Ribeiro, F. Pedregosa, P. van Mulbregt, and SciPy 1.0 Contributors, Nature Methods **17**, 261 (2020).

[43] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," (2017), arXiv:1412.6980 [cs.LG].