Title: NN/QFT correspondence

Speakers: Ro Jefferson

Collection/Series: Theory + AI Workshop: Theoretical Physics for AI

Date: April 11, 2025 - 9:00 AM

URL: https://pirsa.org/25040128

Abstract:

As we've seen at this workshop, exciting progress has recently been made in the study of neural networks by applying ideas and techniques from theoretical physics. In this talk, I will discuss a precise relation between quantum field theory and deep neural networks, the NN/QFT correspondence. In particular, I will go beyond the level of analogy by explicitly constructing the QFT corresponding to a class of networks encompassing both vanilla feedforward and recurrent architectures. The resulting theory closely resembles the well-studied O(N) vector model, in which the variance of the weight initializations plays the role of the 't Hooft coupling. In this framework, the Gaussian process approximation used in machine learning corresponds to a free field theory, and finite-width effects can be computed perturbatively in the ratio of depth to width, T/N. These provide corrections to the correlation length that controls the depth to which information can propagate through the network, and thereby sets the scale at which such networks are trainable by gradient descent. This analysis provides a non-perturbative description of networks at initialization, and opens several interesting avenues to the study of criticality in these models.

NN/QFT Correspondence*

Ro Jefferson (they/them)

Institute for Theoretical Physics, and Department of Information & Computing Sciences, Utrecht University

> Theory + Al Workshop Perimeter Institute, 2025-04-11



k



The state of the art

How deep neural networks work:



"Machine learning has become alchemy" [NeurIPS]

On Friday, someone on another team changed the default rounding mode of some Tensorflow internals (from truncation to "round to even"). Our training broke. Our error rate went from < 25% error to 99.97% error (on a standard 0-1 binary loss).

The state of the art



Need understanding of physical principles, "theory of deep learning"

Physics → Deep Learning



- Shift focus from techniques to boost performance, to simple models to help explain basic phenomena
- Distill fundamental physical principles, simple theory as building blocks to complicated problems
 - \longrightarrow physics-based approach towards a Theory of Deep Learning

k

Applying field-theoretic ideas



Stat phys: scaling limits (Bordelon, Pehlevan), feature learning (Loureiro), spin-glass models (throw a dart)

- RG: Bayesian inference (Berman, Howard), optimal transport (Cotler, Rezchikov), renormalizing Gaussian Processes (Howard, Ringel, Maiti, RJ)
- QFT: effective theory (Roberts, Yaida, Hanin), NN phenomenology (Halverson, Maiti), duality with O(N) models (Grosvenor, RJ)

Infinite-width limit

As $N \to \infty$, statistics Gaussian even if individual z_i^{ℓ} are not (because central limit theorem, provided $\{z_0^{\ell}, z_1^{\ell}, \ldots, z_{N-1}^{\ell}\}$ sequence of i.i.d. random variables with finite variance) \longrightarrow Gaussian Process



Pros: analytically tractable, extremely general (MLP, CNN, RNN, skip connections, etc.)

Cons: infinite-width DNNs effectively shallow, linear models; NTK does not evolve, no feature learning ↔ no interactions



Infinite-width limit

Ŧ

As $N \to \infty$, statistics Gaussian even if individual z_i^{ℓ} are not (because central limit theorem, provided $\{z_0^{\ell}, z_1^{\ell}, \ldots, z_{N-1}^{\ell}\}$ sequence of i.i.d. random variables with finite variance) \longrightarrow Gaussian Process



Pros: analytically tractable, extremely general (MLP, CNN, RNN, skip connections, etc.)

Cons: infinite-width DNNs effectively shallow, linear models; NTK does not evolve, no feature learning ↔ no interactions

Perturbative QFT in a nutshell

- Perturbative QFT: the art of backing away from $N \to \infty$.
- Basic idea: solve complicated (non-Gaussian) theory by perturbing about the free (Gaussian) theory, in terms of some small parameter.
- Physically, turn-on interactions (= finite-width effects) in a controlled manner

Example: quartic interaction

Ŧ

$$p(z) = rac{1}{Z} \exp\left\{-rac{1}{2}z^2 - rac{\lambda}{4!}z^4
ight\} \ , \quad {
m with} \quad \lambda \sim rac{1}{N}$$

Perturbative corrections to correlation functions, e.g.,



Interlude: criticality

Systems at criticality exhibit structure on all scales



Unique trade-off between information transmission and storage



Page 10 of 50

Computation at the edge of chaos

 DNNs are trainable when they lie near criticality chaotic phase: correlations washed-out; hot (exploding gradients) ordered phase: correlations damped; cold (vanishing gradients) critical point: Goldilocks zone between info transmission & storage





1611.01232, 2107.06898

NN/QFT correspondence

Would like to have a bottom-up framework for studying real-world DNNs using techniques/ideas from physics (beyond GP).

Objective: explicitly map DNN to a bona fide QFT



- Starting from NN SDE, marginalize over stochasticity to obtain probability of a particular sequence of network states
- **2** Continuum limit in depth to obtain path integral of O(N) model
- **③** Perturbation theory in L/N at weak 't Hooft coupling σ_w^2
- Compute correlations, study DNN as we do any other field theory

2109.13247, 2505.XXXXX

SDE formulation of RNNs

Statistical field theory approach 1901.10416



$$\mathrm{d}h = f(h,x)\,\mathrm{d}t + g(h,x)\,\mathrm{d}B\;, \qquad h,x,\,\mathrm{d}B \in \mathbb{R}^N$$

 $\begin{array}{ll} \text{Deterministic update:} & f(h,x) = -\gamma h + W \phi(h) + U \varphi(x) + b \\ \\ \text{Stochasticity/noise:} & g(h,x) \in \mathbb{R}^{N \times N} \end{array}$

NN/QFT correspondence

Would like to have a bottom-up framework for studying real-world DNNs using techniques/ideas from physics (beyond GP).

Objective: explicitly map DNN to a bona fide QFT



- Starting from NN SDE, marginalize over stochasticity to obtain probability of a particular sequence of network states
- **2** Continuum limit in depth to obtain path integral of O(N) model
- **③** Perturbation theory in L/N at weak 't Hooft coupling σ_w^2
- Compute correlations, study DNN as we do any other field theory

2109.13247, 2505.XXXXX

SDE formulation of RNNs



$$h_t - h_{t-1} = f(h_{t-1}, x_{t-1})\eta + g_{t-1}\xi_t , \quad \eta \coloneqq \mathrm{d}t , \ \xi_t \coloneqq \frac{\mathrm{d}B_t}{\mathrm{d}t} \mathrm{d}t$$

Assume ξ_t independent; probability of particular path p(t) is

$$p(h(t)) = \int \prod_{t=0}^{T} \mathrm{d}\xi_t \,
ho(\xi_t) \, \delta(h_t - y_t(\xi_t, h_{t-1}))$$

where $\rho(\xi_t)$ is probability density of noise increment ξ_t , and

$$y_t(\xi_t, h_{t-1}) = h_{t-1} + f(h_{t-1}, x_{t-1})\eta + g_{t-1}\xi_t$$

Introduce auxiliary fieldvariable

Express Dirac delta in terms of response field $\tilde{z} = ik \in \mathbb{C}$:

$$\delta(x) = \int_{-\infty}^{\infty} \frac{\mathrm{d}k}{2\pi} \, e^{ikx} = \int_{-i\infty}^{i\infty} \frac{\mathrm{d}\tilde{z}}{2\pi i} \, e^{\tilde{z}x}$$

Probability of path h(t) now an integral over \tilde{z} :

$$p(h(t)) = \prod_{t=0}^{T} \int d\xi_t \,\rho(\xi_t) \int_{-i\infty}^{i\infty} \frac{d\tilde{z}_t}{2\pi i} e^{\tilde{z}_t (h_t - y_t(\xi_t, h_{t-1}))}$$

= $\prod_{t=0}^{T} \int d\xi_t \,\rho(\xi_t) \int_{-i\infty}^{i\infty} \frac{d\tilde{z}_t}{2\pi i} \exp\left[\tilde{z}_t \left(h_t - h_{t-1} - f_{t-1}\eta\right) - \tilde{z}_t g_{t-1}\xi_t\right]$
= $\prod_{t=0}^{T} \int_{-i\infty}^{i\infty} \frac{d\tilde{z}_t}{2\pi i} \exp\left[\tilde{z}_t \left(h_t - h_{t-1} - f_{t-1}\eta\right) + K_{\xi}(-\tilde{z}_t g_{t-1})\right]$

where $K_{\xi}(-\tilde{z}_t g_{t-1}) = \ln \langle e^{-\tilde{z}_t g_{t-1}\xi} \rangle_{\xi} = \ln \int d\xi_t \, \rho(\xi_t) \, e^{-\tilde{z}_t g_{t-1}\xi_t}$ is the cumulant generating function of ξ_t .

Obtain moment generating function

Add source terms $j_t h_t \eta$, integrate over all paths h:

$$\begin{split} Z[j] &= \left\langle \prod_{t=0}^{T} e^{j_t h_t \eta} \right\rangle_h = \prod_{t=0}^{T} \int \mathrm{d}h_t \, p(h(t)) \, e^{j_t h_t \eta} \\ &= \prod_{t=0}^{T} \left\{ \int \mathrm{d}h_t \frac{\mathrm{d}\tilde{z}_t}{2\pi i} \right\} \exp \sum_{t=0}^{T} \left[\tilde{z}_t \left(h_t - h_{t-1} - f_{t-1} \eta \right) + j_t h_t \eta + K_{\xi}(-\tilde{z}_t g_{t-1}) \right] \end{split}$$

Partition function for h: $\partial_{j_t\eta}Z[j]|_{j_t=0} = \langle h_t \rangle_h$

Continuum-limit $\eta \to 0$

Path-integral measures:

$$\lim_{\eta \to 0} \prod_{t=0}^T \int_{-\infty}^{\infty} \mathrm{d}h_t \coloneqq \int \mathcal{D}h \ , \qquad \lim_{\eta \to 0} \prod_{t=0}^T \int_{-i\infty}^{i\infty} \frac{\mathrm{d}\tilde{z}_t}{2\pi i} \eqqcolon \int \mathcal{D}\tilde{z}$$

Fields within the exponential,

$$\lim_{\eta \to 0} \sum_t h_t \eta = \int \mathrm{d}t \, h(t) \,, \qquad \lim_{\eta \to 0} \sum_t \frac{h_t - h_{t-1}}{\eta} \eta = \int \! \mathrm{d}t \, \partial_t h(t)$$

Path integral of continuum field theory:

$$Z[j] = \int \mathcal{D}h \mathcal{D}\tilde{z} \exp\left\{\int dt \left[\tilde{z}(t) \left(\partial_t h(t) - f(t)\right) + j(t)h(t)\right] + K_B(-\tilde{z}g)\right\}$$

where $K_B(-\tilde{z}g) = \ln \int \mathcal{D}\xi \exp\left\{-\int dB \,\tilde{z}(t)g(t)\right\}$



Isolate behaviour/statistics of h by tracing over parameters, obtain partition function for ensemble average $\bar{Z}[j] = \langle Z[j] \rangle_{W,U,b}$

0

Fri 11 Apr 09:39:4

Self-averaging random networks

Consider ensemble of random networks:

$$\begin{aligned} X_{ij} \sim \mathcal{N}(0, \sigma_x^2) , \qquad X = [X_{ij}] \in \{W, U\} ,\\ b_i \sim \mathcal{N}(0, \sigma_b^2) , \qquad b = [b_i] \end{aligned}$$

<u>Self-averaging</u>: instantiations vary, but physical properties given by mean ensemble values:

$$\bar{Z}[j] \coloneqq \langle Z[j] \rangle_{X,b} = \prod_{ij} \int dX_{ij} \int db_i \,\rho(X_{ij}) \,\rho(b_i) \,Z[j]$$
$$=: \int \mathcal{D}X \int \mathcal{D}b \,Z[j]$$

where $ho(X_{ij})=\sqrt{rac{N_x}{2\pi\sigma_x^2}}\,e^{-rac{N}{2}\left(rac{X_{ij}}{\sigma_x}
ight)^2}$

Self-averaging random networks

Plug in f, perform integral over $X \in \{W, U\}$ and b:

$$\bar{Z}[j] = \int \mathcal{D}h \int \mathcal{D}\tilde{z} \, e^{S_0 + S_{
m source} + S_{
m int}}$$

where

$$S_0 = \int \mathrm{d}t \, \tilde{z}(t) (\partial_t + \gamma) \, h(t) + K_B(-\tilde{z}g)$$

$$e^{S_{\text{int}}} = \int \mathcal{D}W \mathcal{D}U \mathcal{D}b \exp\left\{-\int \mathrm{d}t \,\tilde{z}_i(t) \left[W_{ij}\phi(h_j(t)) + U_{ij}\varphi(x_j(t)) + b_i\right]\right\}$$
$$= \left\langle e^{-W_{ij}\int \mathrm{d}t \,\tilde{z}_i\phi(h_j)} \right\rangle_W \left\langle e^{-U_{ij}\int \mathrm{d}t \,\tilde{z}_i\varphi(x_j)} \right\rangle_U \left\langle e^{-b_i\int \mathrm{d}t \,\tilde{z}_i} \right\rangle_b$$

Self-averaging random networks

Each $\langle \ldots \rangle$ a product of i.i.d. Gaussians, e.g.,

$$\left\langle e^{-W_{ij}\int \mathrm{d}t\,\tilde{z}_i\phi(h_j)}\right\rangle_W = \prod_{ij}\int \mathrm{d}W_{ij}\,\rho(W_{ij})\,e^{-W_{ij}\int \mathrm{d}t\,\tilde{z}_i\phi(h_j)}$$
$$= \exp\sum_{ij}\frac{\sigma_w^2}{2N}\int \mathrm{d}t_1\,\mathrm{d}t_2\,\tilde{z}_i(t_1)\,\tilde{z}_i(t_2)\,\phi(h_j(t_1))\,\phi(h_j(t_2))$$

Interaction term then

$$S_{\text{int}} = \frac{1}{2N} \int dt_1 dt_2 \sum_i \tilde{z}_i(t_1) \,\tilde{z}_i(t_2)$$
$$\times \sum_j \left[\sigma_b^2 + \sigma_w^2 \,\phi(h_j(t_1)) \,\phi(h_j(t_2)) + \sigma_u^2 \,\varphi(x_j(t_1)) \,\varphi(x_j(t_2)) \right]$$

N.b., factorized into N independent subsystems \tilde{z}_i^2 coupled to \sum_j

Analogy with O(N) vector model

Motivates introducing auxiliary field variables

$$\mathfrak{W}(t_1, t_2) \coloneqq \sqrt{\frac{\sigma_w^2}{N}} \sum_j \phi(h_j(t_1)) \phi(h_j(t_2)) , \quad \text{sim. } \mathfrak{U}(t_1, t_2)$$

- In O(N), $g_{\mathrm{YM}} = \sqrt{\frac{\lambda}{N}} \implies$ 't Hooft coupling $\lambda \longleftrightarrow \sigma_w^2$
- Convergence requires weak ('t Hooft) coupling, $\sigma_w^2 < \gamma^2$

Enforce constraint via delta functions as before:

$$e^{S_{ ext{int}}} = \int \mathcal{D}\mathfrak{W} \exp\left\{rac{1}{2}\int \mathrm{d}t_1 \,\mathrm{d}t_2 \sum_i ilde{z}_i(t_1) \, ilde{z}_i(t_2) \left[\sigma_b^2 + \sqrt{rac{\sigma_w^2}{N}}\mathfrak{W}(t_1, t_2)
ight]
ight\}
onumber \ imes \delta\left(\mathfrak{W}(t_1, t_2) - \sqrt{rac{\sigma_w^2}{N}} \sum_j \phi(h_j(t_1)) \,\phi(h_j(t_2))
ight)$$

then write delta function as integral over \widetilde{W}

Finally...

Our theory is then

$$ar{Z} = \int \mathcal{D} \mathfrak{X} \, \mathcal{D} \widetilde{X} \, \mathcal{D} h \, \mathcal{D} \widetilde{z} \, e^{S_0 + S_{ ext{int}}}$$

where the quadratic part is

$$S_{0} = \int dt \left[\tilde{z}_{i}(t)(\partial_{t} + \gamma) h_{i}(t) + \frac{\kappa}{2} \tilde{z}_{i}(t) \tilde{z}_{i}(t) \right]$$

+
$$\frac{1}{2} \int dt_{1} dt_{2} \left[\sigma_{b}^{2} + \sqrt{\frac{\sigma_{w}^{2}}{N}} \mathfrak{W}_{0}(t_{1}, t_{2}) + \sqrt{\frac{\sigma_{u}^{2}}{N}} \mathfrak{U}_{0}(t_{1}, t_{2}) \right] \tilde{z}_{i}(t_{1}) \tilde{z}_{i}(t_{2})$$

-
$$\frac{1}{2} \int dt_{1} dt_{2} \left[\widetilde{W}(t_{1}, t_{2}) \mathfrak{W}(t_{1}, t_{2}) + \widetilde{U}(t_{1}, t_{2}) \mathfrak{U}(t_{1}, t_{2}) \right]$$

ß

Perturbation theory

Interaction term intractable, since field h hidden within ϕ :

$$\int \mathrm{d}t_1 \, \mathrm{d}t_2 \sqrt{\frac{\sigma_w^2}{N}} \widetilde{W}(t_1, t_2) \phi(h_i(t_1)) \, \phi(h_i(t_2))$$

Take $\phi(h_i) = \varphi(h_i) = anh(h_i) = h_i - rac{h_i^3}{3} + \mathcal{O}(h_i^5)$

Interaction part then

$$S_{\text{int}} = \frac{1}{2} \int dt_1 dt_2 \left[\sqrt{\frac{\sigma_w^2}{N}} \mathfrak{W}(t_1, t_2) + \sqrt{\frac{\sigma_u^2}{N}} \mathfrak{U}(t_1, t_2) \right] \tilde{z}_i(t_1) \tilde{z}_i(t_2) + \frac{1}{2} \int dt_1 dt_2 \left[\sqrt{\frac{\sigma_w^2}{N}} \widetilde{W}(t_1, t_2) \left(h_i(t_1) h_i(t_2) - \frac{2}{3} h_i(t_1) h_i(t_2)^3 \right) \right] \sqrt{\frac{\sigma_w^2}{N}} \widetilde{U}(t_1, t_2) \left(x_i(t_1) x_i(t_2) - \frac{2}{3} x_i(t_1) x_i(t_2)^3 \right) \right]$$

ß



- h(t) propagating to h(s): $G_{hh}(t-s) = t \longrightarrow s$ • h(t) propagating to $\tilde{z}(s)$: $G_{h\tilde{z}}(t-s) = t \longrightarrow s$
- h(t) propagating to $\tilde{z}(s)$: $G_{h\tilde{z}}(t-s) = t \longrightarrow s$ • $\tilde{z}(t)$ propagating to h(s): $G_{\tilde{z}h}(t-s) = t \longrightarrow s$

•
$$\mathfrak{W}(t_1, t_2)$$
 propagating to $\widetilde{W}(t_3, t_4)$: $G_w = \frac{t_1}{t_2} \longrightarrow \frac{t_3}{t_4}$

 t_1

 t_2

 t_1

 t_2

 t_3 t_4

 t_3 t_4





Infinite $\mathcal{O}(1)$ cacti

First perturbative correction given by infinite sequence of cactus diagrams, each $\mathcal{O}(1)$:



Physical interpretation: statistical fluctuations in ensemble of networks.

 \square

Infinite $\mathcal{O}(T/N)$ mushrooms

Second perturbative correction given by infinite sequence of mushroom diagrams, each $\mathcal{O}(T/N)$:



Physical interpretation: effective interactions due to finite-width effects.

ß

n = 3

Loop-corrected propagator

A more elegant method: diagrammatic recursion relation for linear models



where $X^{(0)}$ and $X^{(1)}$ denote $\mathcal{O}(1)$ and $\mathcal{O}(T/N)$ contributions to X in the expansion

$$X = X^{(0)} + \frac{T}{N}X^{(1)}$$

Ð

Effective correlation length

Infinitesimal perturbations about fixed point:



Examine correlator at small times $|\tau|$:

$$X(\tau) \approx \frac{\kappa}{2} \xi \, e^{-|\tau|/\xi} + \frac{\gamma^2}{2} \xi_0^4 \left(2 + \gamma \frac{T}{N} \xi_0^2 \sigma_w^2\right) \sigma_{b,\text{eff}}^2$$

Identify loop-corrected correlation length:

$$\xi \coloneqq rac{\xi_0}{2} \left[1 + \gamma^2 \xi_0^2 + rac{\gamma T}{8N} (1 + 3\gamma^2 \xi_0^2) \, \xi_0^2 \, \sigma_w^2
ight]$$

ß

Non-linear models \implies 5-pt interaction

Quintic interaction yields much more complicated diagrams, e.g.,



Non-linear models: once more, with feeling!









ß

NN/QFT correspondence (GJ version)

Summary:

- NN/QFT correspondence: explicit duality between systems.
- Systematic computation of finite-width effects, study criticality.
- Remarkable parallels with well-studied O(N) vector model, 't Hooft coupling σ_w^2 ; exploit standard tool set
- Bottom-up approach to physical theory of DNNs.

Open questions:

- Difficult to observe shift in critical point at weak coupling; higher order? Strong coupling? Nonperturbative effects?
- Explore: other observables? RG flow to critical initializations?
- Local rotation symmetry \rightarrow gauge theory, RMT?
- Test it! How big is $\mathcal{O}(1)$? $\mathcal{O}(T/N)$? Loop-corrected ξ vs. empirics?