

How to scale
NNs?

Scaling limits

Conclusion

ResNets of All Shapes and Sizes

Louis-Pierre Chaintron (EPFL)

ProbAI Theory of Scaling Laws Workshop 2026
– University of Warwick, June 23rd

Joint work with Lénaïc Chizat and Javier Maass (EPFL).

Scaling up Neural Networks

How to scale NNs?

ResNets

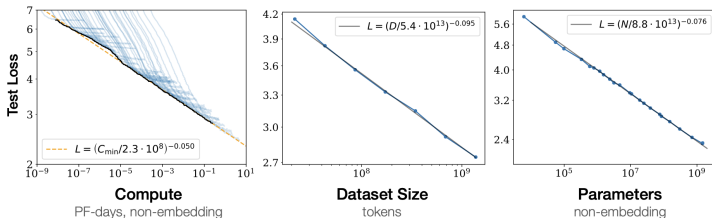
Phase diagram

Scaling limits

Conclusion

Observation:

- ▶ **Scaling up** compute budget of training NNs improves performance.
- ▶ Several ways : more data, longer training, **bigger models**.



Performance vs compute [Kaplan et al'20]
 N parameters, dataset size D , compute C

- ↪ **Classify** large-scale limits wrt hyper-parameter (HP) scalings.
- ⇒ **Improve choices of HP.**

Residual Neural Networks

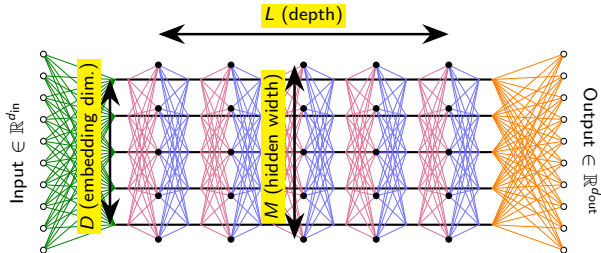
How to scale
NNs?

ResNets

Phase diagram

Scaling limits

Conclusion



$$\left\{ \begin{array}{l} \underbrace{h^0(x)}_{\text{input}} = W_{\text{in}} x \in \mathbb{R}^D, \\ h^\ell(x) = h^{\ell-1}(x) + \frac{1}{LM} \sum_{j=1}^M v^{j,\ell} \underbrace{\rho(u^{j,\ell} \cdot h^{\ell-1}(x))}_{\text{pre-activation}}, \quad \ell \in [1:L], \\ y(x) = W_{\text{out}}^T h^L(x) \in \mathbb{R}^{d_{\text{out}}}. \end{array} \right.$$

$\rho: \mathbb{R} \rightarrow \mathbb{R}$ (activation)

Parameters θ :

$$W_{\text{in}} \in \mathbb{R}^{d_{\text{in}} \times D}, \quad (u^{j,\ell}, v^{j,\ell})_{j,\ell} \in (\mathbb{R}^D \times \mathbb{R}^D)^{M \times L}, \quad W_{\text{out}} \in \mathbb{R}^{d_{\text{out}} \times D}.$$

Training dynamics

How to scale NNs?

ResNets

Phase diagram

Scaling limits

Conclusion

- ▶ Training set $(x_i)_{1 \leq i \leq n}$, loss function:

$$\mathcal{L}(\theta) = \frac{1}{n} \sum_{i=1}^n \text{Loss}_i(y(x_i)).$$

- ▶ Gradient descent:

$$\begin{cases} W_0^{\text{in}}[i, j] \stackrel{iid}{\sim} \mathcal{SG}(0, \sigma_{\text{in}}^2), \\ u_0^{j, \ell}[i] \stackrel{iid}{\sim} \mathcal{SG}(0, \sigma_u^2), \\ v_0^{j, \ell}[i] \stackrel{iid}{\sim} \mathcal{SG}(0, \sigma_v^2), \\ W_0^{\text{out}}[i, j] \stackrel{iid}{\sim} \mathcal{SG}(0, \sigma_{\text{out}}^2), \end{cases} \quad \begin{cases} W_{k+1}^{\text{in}} = W_k^{\text{in}} - \eta_{\text{in}} \nabla_{W^{\text{in}}} \mathcal{L}(\theta_k), \\ u_{k+1}^{j, \ell} = u_k^{j, \ell} - \eta_u \nabla_{u^{j, \ell}} \mathcal{L}(\theta_k), \\ v_{k+1}^{j, \ell} = v_k^{j, \ell} - \eta_v \nabla_{v^{j, \ell}} \mathcal{L}(\theta_k), \\ W_{k+1}^{\text{out}} = W_k^{\text{out}} - \eta_{\text{out}} \nabla_{W^{\text{out}}} \mathcal{L}(\theta_k). \end{cases}$$

Given $(\sigma_{\text{in}}, \sigma_u, \sigma_v, \sigma_{\text{out}})$, we scale $(\eta_{\text{in}}, \eta_u, \eta_v, \eta_{\text{out}})$ so that:

- ▶ Loss decay due to $W_{\text{in}}, W_{\text{out}}$ is $\Theta(1)$.
- ▶ Loss decay due to $u^{j, \ell}, v^{j, \ell}$ is $\Theta(1/(LM))$.

Pre-activation in $\Theta(1)$: we choose $\sigma_{\text{in}} \sim 1$ and $\sigma_u = \sigma_v \sim 1/\sqrt{D}$.

Phase diagram (assuming $ML \gg D \gg 1$)

How to scale NNs?

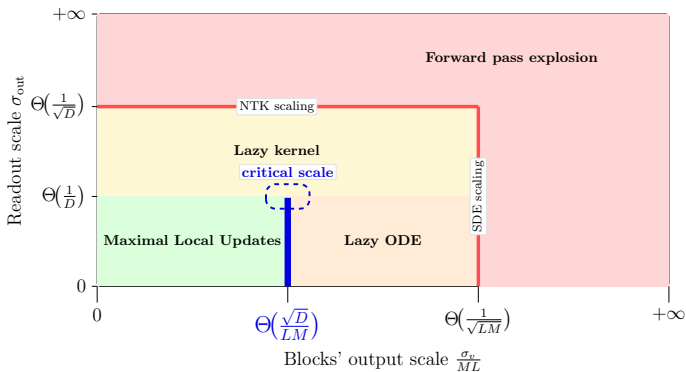
ResNets

Phase diagram

Scaling limits

Conclusion

- ▶ **Lazy kernel:** (pre)activations update in $\mathcal{O}(1)$ at any step.
- ▶ **Feature learning:** (pre)activations update in $\Theta(1)$ in the first GD step.
- ▶ **Maximal Local Update (MLU):** contribution in $\Theta(1)$ of local weights to update.
- ▶ **Lazy ODE:** feature learning but not local – local contributions in $\mathcal{O}(1)$.



[Chizat et al.'18], [Yang et al.'21], [Dey et al.'25], [Chizat'25]...

Application: HP transfer [Yang et al'22]

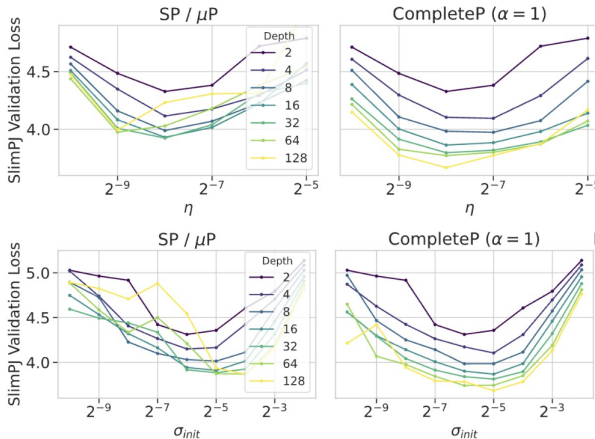
How to scale NNs?

ResNets

Phase diagram

Scaling limits

Conclusion



From “CompleteP” paper [Cerebras AI et al., '25]. Transformer of varying depth trained on 300M tokens – here CompleteP = critical scale.

Large-scale limit in the MLU regime

How to scale NNs?

ResNets

Phase diagram

Scaling limits

Conclusion

	Tensor Program/DMFT	Neural Mean ODE
Approach	$M \propto D \rightarrow \infty$ then $L \rightarrow \infty$	$M, L \rightarrow \infty$ with D fixed
Other limitations	qualitative, 2nd step heuristic	loose upper bound
Main refs	[Yang et al., Bordelon et al.,...]	[Lu et al., Ding et al.,...]

Nb of params.	8B	70B	405B
L	32	80	126
D	4,096	8,192	16,384
M	14,336	28,672	53,248
M/D	3.5	3.5	3.25
ML/D	112	280	410

Shape hyperparameters of Llama 3.1 (MLP blocks)

↔ Unified theory for joint limits $L, M, D \rightarrow +\infty$? Quantitative rigorous rates?

How to scale
NNs?

Scaling limits

Mean ODE

Large D

DMFT

Conclusion

Scaling limits

Limit theorem

How to scale

NNs?

Scaling limits

Mean ODE

Large D

DMFT

Conclusion

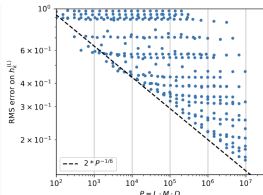
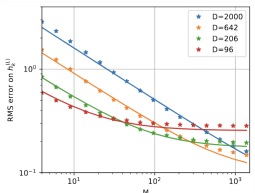
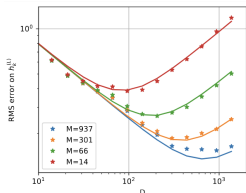
Theorem (C., Chizat, Maass'26)

\exists (limit dynamics) : $\forall k \geq 1, \exists c_k, c'_k > 0 : \forall \delta \in (0, 1],$

$$(\text{error})_k \leq c_k \left[\frac{1}{L} + \sqrt{\frac{D}{LM}} + \frac{1 + \log(1/\delta)}{\sqrt{D}} \right] \quad \text{with } \text{proba} \geq 1 - \delta,$$

provided $RHS \leq c'_k$ – for clipped GD.

\hookrightarrow For a budget $P = LMD$, we get $(\text{error}) = O(P^{-1/6})$.



Error in embedding space for fitted coefficients at $k = 15$.

Backward pass

How to scale

NNs?

Scaling limits

Mean ODE

Large D

DMFT

Conclusion

Training dynamics:

$$\left\{ \begin{array}{l} h_k^0 = W_{\text{in}} x, \\ h_k^\ell = h_k^{\ell-1} + \frac{1}{LM} \sum_{j=1}^M \rho(\langle u_k^{j,\ell}, h_k^{\ell-1} \rangle_{\overline{D}}) v_k^{j,\ell}, \\ f_k(x) = W_{\text{out}}^\top h_k^\ell, \\ b_k^{\ell-1} = b_k^\ell + \frac{1}{LM} \sum_{j=1}^M \rho'(\langle u_k^{j,\ell}, h_k^{\ell-1} \rangle_{\overline{D}}) \langle v_k^{j,\ell}, b_k^{\ell-1} \rangle_{\overline{D}} u_k^{j,\ell}, \\ b_k^\ell = W_{\text{out}}^\top \nabla \mathcal{L}(f_k(x)), \\ u_{k+1}^{j,\ell} = u_k^{j,\ell} - \eta_u \rho'(\langle u_k^{j,\ell}, h_k^{\ell-1} \rangle_{\overline{D}}) \langle v_k^{j,\ell}, b_k^{\ell-1} \rangle_{\overline{D}} h_k^{\ell-1}, \\ v_{k+1}^{j,\ell} = v_k^{j,\ell} - \eta_v \rho(\langle u_k^{j,\ell}, h_k^{\ell-1} \rangle_{\overline{D}}) b_k^{\ell-1}. \end{array} \right.$$

Normalized backward pass:

$$b_k^\ell := D W_{\text{out}} \left[\frac{\partial h_k^\ell}{\partial h_k^\ell} \right]^\top \nabla \mathcal{L}(f_k(x)), \quad \begin{cases} u_0^{j,\ell} \sim \mathcal{SG}(0, D\sigma_u^2), \\ v_0^{j,\ell} \sim \mathcal{SG}(0, D\sigma_v^2). \end{cases}$$

Neural Mean ODE

How to scale NNs?

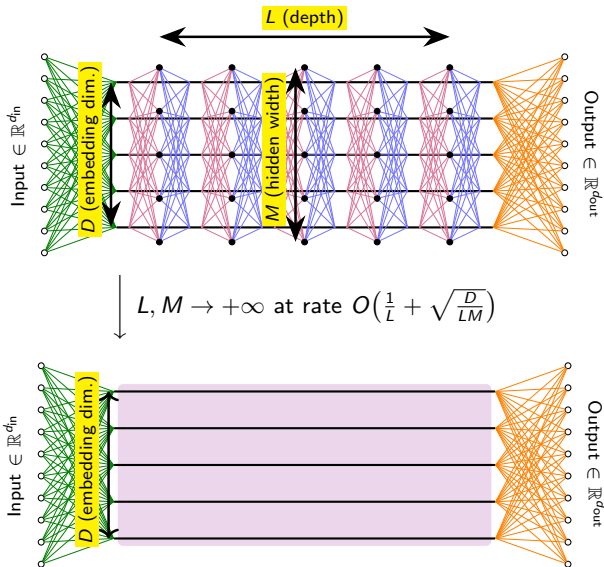
Scaling limits

Mean ODE

Large D

DMFT

Conclusion



Neural Mean ODE

How to scale NNs?

Scaling limits

Mean ODE

Large D

DMFT

Conclusion

Limit $L, M \rightarrow +\infty$ dynamics:

$$\begin{cases} \partial_s h_k(s) = \mathbb{E}[V_k(s) \rho(\langle U_k(s), h_k(s) \rangle_{\overline{D}})], \\ \partial_s b_k(s) = -\mathbb{E}[V_k(s) \rho'(\langle U_k(s), h_k(s) \rangle_{\overline{D}}) \langle V_k(s), b_k(s) \rangle_{\overline{D}} U_k(s)], \\ U_{k+1}(s) = U_k(s) - \eta_u \rho'(\langle U_k(s), h_k(s) \rangle_{\overline{D}}) \langle V_k(s), b_k(s) \rangle_{\overline{D}} h_k(s), \\ V_{k+1}(s) = V_k(s) - \eta_v \rho(\langle U_k(s), h_k(s) \rangle_{\overline{D}}) b_k(s). \end{cases}$$

Theorem (Chizat'25)

For *Clipped GD* in the *MLU* regime, if $\log L \vee D \lesssim M$, then with high probability

$$\max_{0 \leq \ell \leq L} |h_k^\ell(x) - h_k(\ell/L, x)| = O\left(\frac{1}{L} + \sqrt{\frac{D}{LM}}\right).$$

Experiment: narrow vs deep ResNet

How to scale NNs?

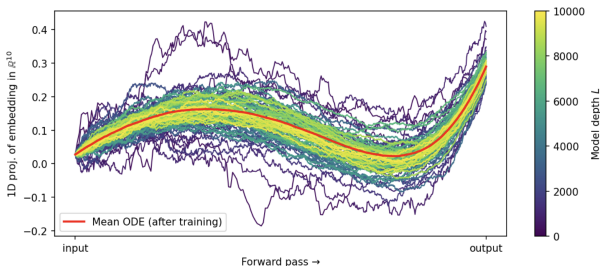
Scaling limits

Mean ODE

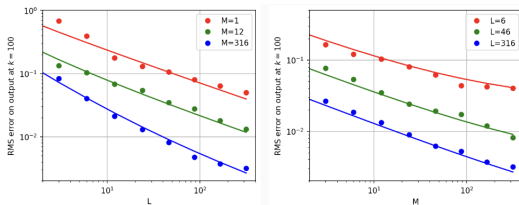
Large D

DMFT

Conclusion



1D projection of the forward pass on one sample after 50 GD steps
– GD for the square-loss on $n = 10$ training samples with $D = 10$, $M = 1$.



Experiment vs theory (with fitted coefficients) after $k = 100$ GD steps.

Large D limit

How to scale NNs?

Scaling limits

Mean ODE

Large D

DMFT

Conclusion

Neural mean ODE in \mathbb{R}^D :

$$\begin{cases} \partial_s h_k(s) = \mathbb{E}[V_k(s) \rho(\langle U_k(s), h_k(s) \rangle_{\overline{D}})], \\ \partial_s b_k(s) = -\mathbb{E}[V_k(s) \rho'(\langle U_k(s), h_k(s) \rangle_{\overline{D}}) \langle V_k(s), b_k(s) \rangle_{\overline{D}} U_k(s)], \\ U_{k+1}(s) = U_k(s) - \eta_u \rho'(\langle U_k(s), h_k(s) \rangle_{\overline{D}}) \langle V_k(s), b_k(s) \rangle_{\overline{D}} h_k(s), \\ V_{k+1}(s) = V_k(s) - \eta_v \rho(\langle U_k(s), h_k(s) \rangle_{\overline{D}}) b_k(s). \end{cases}$$

Initialisation:

$$U_k(s) = \underbrace{U_0(s)} + \underbrace{\Delta U_k(s)} = \sqrt{D} U + \Delta U_k(s),$$

$$\mathcal{SG}(0, D\sigma_u^2) \quad O(1)$$

so that

$$\langle U_k(s), h_k(s) \rangle_{\overline{D}} = \frac{1}{\sqrt{D}} \sum_{d=1}^D U^d h_k^d(s) + \frac{1}{D} \sum_{d=1}^D \Delta U_k^d(s) h_k^d(s),$$

for $(U^d)_{d \geq 1}$ i.i.d. $\sim \mathbf{1}$.

Skeleton maps

How to scale NNs?

Scaling limits

Mean ODE

Large D

DMFT

Conclusion

CLT-type sums

$$(\mathbf{S}^{\mathbf{h}^D}(s), \mathbf{S}^{\mathbf{b}^D}(s)) := \left(\left(\frac{1}{\sqrt{D}} \sum_{d=1}^D U^d h_i^d(s), \frac{1}{\sqrt{D}} \sum_{d=1}^D V^d b_i^d(s) \right) \right)_{0 \leq i \leq K-1}$$

Skeleton maps

$$(\Delta U_k^{(D)}(s), \Delta V_k^{(D)}(s)) = (f_k^{\mathbf{h}^D, \mathbf{b}^D}, g_k^{\mathbf{h}^D, \mathbf{b}^D})(s, \mathbf{S}^{\mathbf{h}^D}(s), \mathbf{S}^{\mathbf{b}^D}(s))$$

Definition

Functions $\mathbf{f}^{\mathbf{h}^D, \mathbf{b}^D}, \mathbf{g}^{\mathbf{h}^D, \mathbf{b}^D} : [0, 1] \times \mathbb{R}^K \times \mathbb{R}^K \rightarrow (\mathbb{R}^D)^K$ such that

$$\begin{cases} f_{k+1}^{\mathbf{h}^D, \mathbf{b}^D} = f_k^{\mathbf{h}^D, \mathbf{b}^D} - \eta_u \rho'(z_k^h + \langle h_k^D, f_k^{\mathbf{h}^D, \mathbf{b}^D} \rangle_{\overline{D}}) [z_k^b + \langle b_k^D, g_k^{\mathbf{h}^D, \mathbf{b}^D} \rangle_{\overline{D}}] h_k^D, \\ g_{k+1}^{\mathbf{h}^D, \mathbf{b}^D} = g_k^{\mathbf{h}^D, \mathbf{b}^D} - \eta_v \rho(z_k^h + \langle h_k^D, f_k^{\mathbf{h}^D, \mathbf{b}^D} \rangle_{\overline{D}}) b_k^D, \end{cases}$$

evaluating at $(s, \mathbf{z}^h, \mathbf{z}^b)$.

→ As $D \rightarrow +\infty$, coordinates become **independent**.

Limit skeleton

How to scale

NNs?

Scaling limits

Mean ODE

Large D

DMFT

Conclusion

Definition (Mean-field Skeleton)

Random maps $\mathbf{F}^{\mathbf{H},\mathbf{B}}, \mathbf{G}^{\mathbf{H},\mathbf{B}} : [0, 1] \times \mathbb{R}^K \times \mathbb{R}^K \rightarrow \mathbb{R}^K$ given by

$$\begin{cases} F_{k+1}^{\mathbf{H},\mathbf{B}} = F_k^{\mathbf{H},\mathbf{B}} - \eta_u \rho'(z_k^h + \mathbb{E}[H_k F_k^{\mathbf{H},\mathbf{B}}]) [z_k^b + \mathbb{E}[B_k G_k^{\mathbf{H},\mathbf{B}}]] H_k, \\ G_{k+1}^{\mathbf{H},\mathbf{B}} = G_k^{\mathbf{H},\mathbf{B}} - \eta_v \rho(z_k^h + \mathbb{E}[H_k F_k^{\mathbf{H},\mathbf{B}}]) B_k, \end{cases}$$

The finite- D dynamics, $1 \leq d \leq D$,

$$\begin{aligned} \partial_s h_k^d(s) &= \mathbb{E}[V_k^d(s) \rho(\langle U_k(s), h_k(s) \rangle_{\overline{D}})] \\ &= \mathbb{E}[(\sqrt{D} V^d + (g_k^{\mathbf{h}^D, \mathbf{b}^D})^d) \rho(\langle \sqrt{D} U_k + f_k^{\mathbf{h}^D, \mathbf{b}^D}, h_k(s) \rangle_{\overline{D}})], \end{aligned}$$

becomes, as $D \rightarrow +\infty$,

$$\begin{aligned} \partial_s H_k(s) &= \sigma_v^2 \mathbb{E}[\rho'(P_k(s)) \mathbb{E}[H_k(s) \nabla_{z^b} F_k^{\mathbf{H},\mathbf{B}}(s, \mathbf{Z}^h, \mathbf{Z}^b) | \mathbf{Z}^h, \mathbf{Z}^b]] \cdot \mathbf{B}_{\wedge k-1}(s) \\ &\quad + \mathbb{E}[\rho(P_k(s)) G_k^{\mathbf{H},\mathbf{B}}(s, \mathbf{Z}^h, \mathbf{Z}^b) | W_{\text{in}}, W_{\text{out}}], \end{aligned}$$

$$P_k(s) := Z_k^h(s) + \mathbb{E}[H_k(s) F_k^{\mathbf{H},\mathbf{B}}(s, \mathbf{Z}^h, \mathbf{Z}^b) | \mathbf{Z}^h, \mathbf{Z}^b].$$

Linear setting

How to scale NNs?

Scaling limits

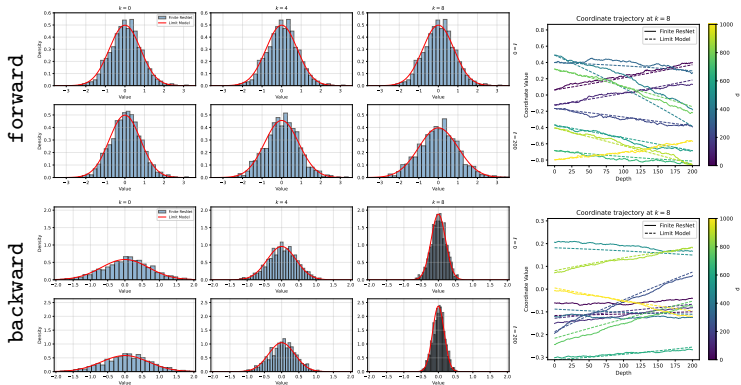
Mean ODE

Large D

DMFT

Conclusion

In the linear case $\rho(z) = z$, the limit is explicit.



→ This matches the behaviour of finite ResNets.

DMFT Closure

How to scale NNs?

Scaling limits

Mean ODE

Large D

DMFT

Conclusion

We only need evaluation at $(\mathbf{Z}^h(s), \mathbf{Z}^b(s))$. Setting

$$\begin{cases} (F_k, G_k) := (F_k^{\mathbf{H}, \mathbf{B}}(\mathbf{Z}^h, \mathbf{Z}^b), G_k^{\mathbf{H}, \mathbf{B}}(\mathbf{Z}^h, \mathbf{Z}^b)), \\ (\nabla_{\mathbf{z}^j} F_k, \nabla_{\mathbf{z}^j} G_k) := (\nabla_{\mathbf{z}^j} F_k^{\mathbf{H}, \mathbf{B}}(\mathbf{Z}^h, \mathbf{Z}^b), \nabla_{\mathbf{z}^j} G_k^{\mathbf{H}, \mathbf{B}}(\mathbf{Z}^h, \mathbf{Z}^b)), \quad j \in \{h, b\}, \end{cases}$$

we have the update rules

$$\begin{cases} F_{k+1} = F_k - \eta_u \rho'(P_k) Q_k H_k, \\ G_{k+1} = G_k - \eta_v \rho(P_k) B_k, \end{cases}$$

$$\begin{aligned} \nabla_{\mathbf{z}^j} F_{k+1} = & \left(\nabla_{\mathbf{z}^j} F_k \right) - \eta_u \left(\rho''(P_k) Q_k \left(\mathbb{E} [H_k \nabla_{\mathbf{z}^j} F_k | \mathbf{z}^h, \mathbf{z}^b] \right) \right. \\ & \left. + \rho'(P_k) \left(\mathbb{E} [B_k \nabla_{\mathbf{z}^j} G_k | \mathbf{z}^h, \mathbf{z}^b] \right) \right) H_k, \end{aligned}$$

$$\nabla_{\mathbf{z}^j} G_{k+1} = \left(\nabla_{\mathbf{z}^j} G_k \right) - \eta_v \rho'(P_k) \left(\mathbb{E} [H_k \nabla_{\mathbf{z}^j} F_k | \mathbf{z}^h, \mathbf{z}^b] \right) B_k, \quad \text{for } j \in \{h, b\}.$$

↪ **Closure** in the variables $(W_{\text{in}}, W_{\text{out}}, \mathbf{H}, \mathbf{B}, \mathbf{F}, \mathbf{G}, \nabla \mathbf{F}, \nabla \mathbf{G})$.

↪ Similar to **linear response functions** in DMFT [Bordelon et al. '22, Montanari et al.'25,...].

How to scale

NNs?

Scaling limits

Mean ODE

Large D

DMFT

Conclusion

Theorem (C., Chizat, Maass'26)

The limit system admits a *unique solution*
 $(\mathbf{H}, \mathbf{B}) \in L^2(\Omega, C([0, 1], \mathbb{R}^{2K}))$.

Theorem (C., Chizat, Maass'26)

Let h_k^D and H_k^D be *coupled* through the same $(W_{\text{in}}^D, W_{\text{out}}^D)$.

$\exists C_k > 0 : \forall D \geq 1, \forall \delta \in (e^{-D}, 1)$,

$$\sup_{s \in [0, 1]} \|h_k^D(s) - H_k^D(s)\|_{\bar{D}} \leq C_k \frac{1 + \log(1/\delta)}{\sqrt{D}},$$

with *proba* $\geq 1 - \delta$.

If furthermore $(W_{\text{in}}^D, W_{\text{out}}^D)$ is a.s. bounded, then this holds for $\delta \in (0, 1)$, yielding *L^2 -convergence*.

Cavity method

How to scale NNs?

Scaling limits

Mean ODE

Large D

DMFT

Conclusion

Taylor expansion:

$$\begin{aligned} & \mathbb{E}[\rho(\langle U_k^D, h_k^D \rangle_{\overline{D}}) \sqrt{D} V^d | \mathbf{h}, \mathbf{b}] \\ &= \mathbb{E}[\rho(\mathbf{S}_k^{\mathbf{h}^D} + \langle f_k^{\mathbf{h}^D, \mathbf{b}^D}(\mathbf{S}^{\mathbf{h}^D}, \mathbf{S}^{\mathbf{b}^D}), h_k^D \rangle_{\overline{D}}) \sqrt{D} V^d | \mathbf{h}, \mathbf{b}] \\ &\simeq \mathbb{E}[\rho(\mathbf{S}_k^{\mathbf{h}^D} + \langle f_k^{\mathbf{h}^D, \mathbf{b}^D}(\mathbf{S}^{\mathbf{h}^D}, \overline{\mathbf{S}}^{\mathbf{b}^D}), h_k^D \rangle_{\overline{D}}) \sqrt{D} V^d | \mathbf{h}, \mathbf{b}] \\ &+ \mathbb{E}[\rho'(\mathbf{S}_k^{\mathbf{h}^D} + \langle f_k^{\mathbf{h}^D, \mathbf{b}^D}(\mathbf{S}^{\mathbf{h}^D}, \overline{\mathbf{S}}^{\mathbf{b}^D}), h_k^D \rangle_{\overline{D}}) \langle \nabla_{z^{\mathbf{b}}} f_k^{\mathbf{h}^D, \mathbf{b}^D}(\mathbf{S}^{\mathbf{h}^D}, \overline{\mathbf{S}}^{\mathbf{b}^D}) [V^d \mathbf{b}^d], h_k^D \rangle_{\overline{D}} V^d] \\ &+ O(D^{-1/2}), \end{aligned}$$

where $\overline{\mathbf{S}}^{\mathbf{b}^D} := \mathbf{S}^{\mathbf{b}^D} - V^d \mathbf{b}^d / \sqrt{D}$, using

$$(\mathbf{S}^{\mathbf{h}^D}, \overline{\mathbf{S}}^{\mathbf{b}^D}) \perp\!\!\!\perp V^d \quad | \quad (\mathbf{h}, \mathbf{b}).$$

As $D \rightarrow +\infty$, we get

$$\sigma_v^2 \mathbb{E}[\rho'(P_k)] \mathbb{E}[H_k \nabla_{z^{\mathbf{b}}} F_k^{\mathbf{H}, \mathbf{B}}(\mathbf{Z}^{\mathbf{H}}, \mathbf{Z}^{\mathbf{B}}) | \mathbf{Z}^{\mathbf{H}}, \mathbf{Z}^{\mathbf{B}}] | \mathbf{H}, \mathbf{B}] \mathbf{B}.$$

+ Quantitative CLT.

Conclusion

How to scale

NNs?

Scaling limits

Conclusion

Highlights: a theory of ResNets in the **large-scale** limit

- ▶ **Phase diagram** of ResNets.
- ▶ Infinite depth ($L \rightarrow +\infty$) implies infinite width ($M \rightarrow +\infty$).
- ▶ High-dimensional ($D \rightarrow +\infty$) limit of Mean ODE.
- ▶ Joint limit at rate $O\left(\frac{1}{L} + \sqrt{\frac{D}{LM}} + \frac{1}{\sqrt{D}}\right)$ – numerically sharp.

Perspectives: complete phase diagram, learning properties of the limit, tighten connections to practice. . .

Articles:

- ▶ Chizat (2025). The Hidden Width of Deep ResNets: Tight Error Bounds and Phase Diagrams <arXiv:2509.10167>.
- ▶ C., Chizat, Maass (2026). ResNets of All Shapes and Sizes: Convergence of Training Dynamics in the Large-scale Limit <arXiv:2603.18168>.