
Universal Sharpness Dynamics in Neural Network Training: Fixed Point Analysis, Edge of Stability, and Route to Chaos

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In gradient descent dynamics of neural networks, the top eigenvalue of the loss
2 Hessian (sharpness) displays a variety of robust phenomena throughout training.
3 This includes early time regimes where the sharpness may decrease during early
4 periods of training (sharpness reduction), and later time behavior such as pro-
5 gressive sharpening and edge of stability. We demonstrate that a simple 2-layer
6 linear network (UV model) trained on a single training example exhibits all of the
7 essential sharpness phenomenology observed in real-world scenarios. By analyzing
8 the structure of dynamical fixed points in function space and the vector field of
9 function updates, we uncover the underlying mechanisms behind these sharpness
10 trends. Our analysis reveals (i) the mechanism behind early sharpness reduction
11 and progressive sharpening, (ii) the required conditions for edge of stability, (iii)
12 the crucial role of initialization and parameterization, and (iv) a period-doubling
13 route to chaos on the edge of stability manifold as learning rate is increased. Finally,
14 we demonstrate that various predictions from this simplified model generalize to
15 real-world scenarios and discuss its limitations.

16 1 Introduction

17 Over the last several years, it has been observed that the training dynamics of neural networks
18 exhibits a rich and robust set of unexpected phenomena, stemming from the non-convexity of the loss
19 landscape. These phenomena not only challenge our existing understanding of loss landscapes but
20 also open avenues for significantly enhancing model performance through improved optimization
21 techniques. In particular, the unexpected and robust phenomenology is mainly associated with the
22 evolution of the Hessian of the loss function, which provides a measure of the local curvature of the
23 loss landscape and plays an important role in understanding generalization performance [18, 11, 16].

24 On the one hand, it has been observed that at late training times, gradient descent (GD) typically
25 exhibits “progressive sharpening,” where the top eigenvalue of the loss Hessian λ^H , referred to as
26 the sharpness, gradually increases with time, until it reaches roughly $2/\eta$, where η is the learning
27 rate. Once the sharpness reaches roughly $2/\eta$, it stops increasing and typically oscillates near $2/\eta$, a
28 late-time training phenomenon referred to as the “edge of stability (EoS)” [7]. On the other hand,
29 during early training, a decrease in sharpness is observed —referred to as “sharpness reduction” [17]
30 —before hitting a temporary plateau.

31 For large enough learning rates, training temporarily destabilizes early on, and the network “catapults”
32 out of its local basin, leading to a temporary sudden increase in the loss in the first few steps, before
33 eventually settling down in a flatter region of the loss landscape characterized by lower sharpness [22].
34 Similar to the loss, sharpness may also spike within the first few steps of training and quickly decrease

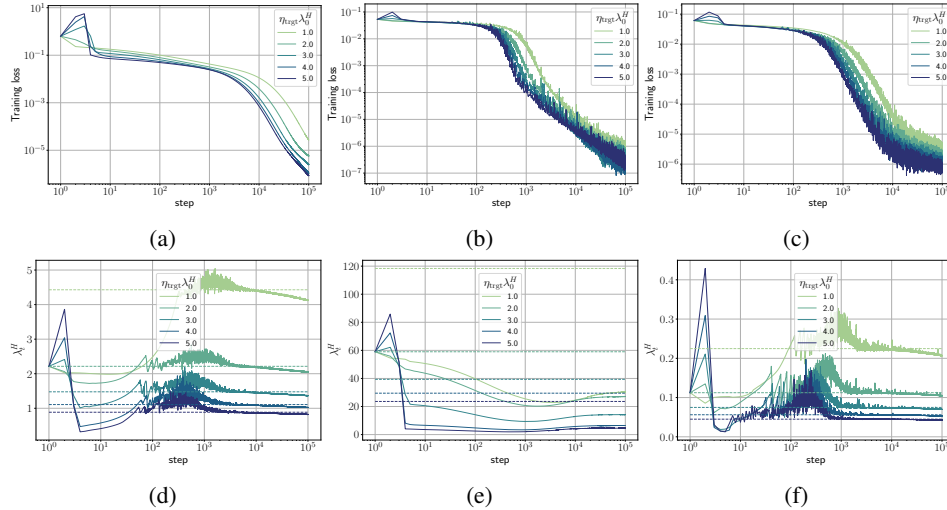


Figure 1: Training loss and sharpness trajectories of ReLU FCNs trained on a 5k subset of CIFAR-10 examples using MSE loss and GD: (a, d) SP with $\sigma_w^2 = 0.5$, (b, e) SP with $\sigma_w^2 = 2.0$, (c, f) μP with $\sigma_w^2 = 2.0$. The dashed lines in the sharpness figures show the $2/\eta$ threshold.

35 (sharpness catapult). A rich phase diagram as a function of network depth, width and learning rate
 36 summarizes the early training dynamics [17].

37 The discovery of these intriguing sharpness phenomena has attracted significant attention, with an
 38 emphasis on various toy models that exhibit similar phenomenology. Yet, the specific conditions and
 39 reasons why these phenomena occur still remain elusive. In this paper, we analyze a simple toy model,
 40 a 2-layer linear network trained on one example, referred to as the UV model. We show that all of the
 41 phenomena described above can be observed in the UV model for appropriate choices of learning rate,
 42 initialization, parameterization, and choice of training example. Through this exploration, we
 43 provide novel insights into the mechanisms at play and offer predictions that we validate in realistic
 44 architectures with both real and synthetic datasets.

45 **Our Contributions.** We revisit the four training regimes identified by Ref. [17] (early time transient,
 46 intermediate saturation, progressive sharpening, and late time EoS) in Section 3, focusing on the
 47 crucial role of initializations and parameterizations. Our findings reveal that models in Standard
 48 Parameterization (SP) with large initializations do not exhibit EoS, even at late training times.
 49 Moreover, we show that models in Maximal Update Parameterization (μP) [35] do not experience an
 50 early sharpness reduction. This result also holds for models in SP with small initializations.

51 We show the UV model exhibits all four training regimes and also captures the effect of initializations
 52 and parameterization discussed above. Through fixed-point analysis of the UV model in the function
 53 space, we analyze the origins of the various dynamical phenomena exhibited by the sharpness.
 54 Specifically, we demonstrate in Sections 4 and 5: (i) the emergence of various sharpness phenomena
 55 arising from the stability and position of the dynamical fixed points, (ii) a critical learning rate η_c ,
 56 above which the model exhibits EoS on a sub-quadratic manifold, and (iii) a period-doubling route
 57 to chaos of sharpness fluctuations as learning rate is increased in the EoS regime. In Appendix A,
 58 we verify various non-trivial predictions from the UV model in realistic architectures with real and
 59 synthetic datasets. Our findings reveal: (i) a sharpness-weight norm correlation before the training
 60 enters the EoS regime, (ii) a phase diagram of EoS, revealing initializations and parameterizations
 61 that do not exhibit EoS, and (iii) a period-doubling route to chaos in real architectures trained on
 62 synthetic datasets, while those trained on real datasets exhibit long-range correlations at the EoS,
 63 with a remnant of the period doubling route to chaos.

64 Given that our analysis spans the entire training trajectory, it relates to numerous studies. Hence, we
 65 defer a comprehensive discussion of related works to Appendix B.

66 2 Notations and Preliminaries

67 This section describes the fundamental concepts and notations that form the basis of our analysis.

68 **Dynamical Systems and Fixed Points:** Consider a discrete dynamical system described by $\theta_{t+1} =$
69 $M(\theta_t)$. A fixed point θ^* of the dynamics satisfies $M(\theta^*) = \theta^*$. The linear stability of a fixed point
70 θ^* is determined by analyzing the eigenvalues $\{\lambda_i^{J^*}\}$ of the Jacobian $J_M(\theta^*) := \nabla_{\theta} M(\theta) |_{\theta=\theta^*}$.
71 An eigendirection $u_i^{J^*}$ of a fixed point θ^* is stable if $|\lambda_i^{J^*}| < 1$ and unstable if $|\lambda_i^{J^*}| > 1$ [25]. The
72 dynamics is captured by the vector field of updates $G(\theta) := M(\theta) - \theta$. The corresponding unit vector
73 is denoted $\hat{G}(\theta) := G(\theta)/\|G(\theta)\|$. Nullclines refer to curves where one of the variables, θ_i , remains
74 invariant, i.e., $\theta_{i;t} = M_i(\theta_t)$.

75 **Parameterizations in Neural Networks:** Sharpness phenomena in neural networks are intrinsically
76 tied to network parameterization. Standard Parameterization (SP) [29] and Neural Tangent Parame-
77 terization (NTP) [14] are two commonly used parameterizations, which converge to kernel methods
78 at infinite width. Ref. [35] proposed Maximal update Parameterization (μ P), which allows for feature
79 learning at infinite width. For implementation details, see Appendix C.2.1.

80 **UV Model:** The UV model refers to a 2-layer linear network $f : \mathbb{R}^d \rightarrow \mathbb{R}$ trained on a single example.
81 We parameterize f as $f(x; \theta) = \frac{1}{\sqrt{n^{1-p}}} v^T U x$, where $x \in \mathbb{R}^d$ is the input, n is the network width,
82 and $v \in \mathbb{R}^n$, $U \in \mathbb{R}^{n \times d}$ are trainable parameters, with each component drawn i.i.d. at initialization
83 from a normal distribution $\mathcal{N}(0, \sigma_w^2/n^p)$. Here, $p \in [0, 1]$ is a parameter that interpolates between
84 NTP and μ P, and $n_{\text{eff}} := n^{1-p}$ is referred to as the effective width. We consider the network trained
85 on a single training example (x, y) using MSE loss $\ell(f(x; \theta), y) = \frac{1}{2} (f(x; \theta) - y)^2$.

86 3 Review of The Four Regimes of Training

87 Typical training trajectories of neural networks can be categorized into four training regimes [17], as
88 shown in Figure 1(a, d):

89 (T1) *Early time transient:* This corresponds to the first few steps of training. At small learning rates
90 ($\eta < \eta_{\text{loss}}$), loss and sharpness decrease monotonically. At larger learning rates ($\eta > \eta_{\text{loss}}$), training
91 catapults out of the initial basin, temporarily increasing the loss, and finally converges to a flatter
92 region [22]. By the end of this regime, sharpness has decreased from initialization for all learning
93 rates, and more substantially at larger learning rates.

94 (T2) *Intermediate saturation:* Following the initial transient regime, sharpness approximately plateaus
95 before gradually increasing.

96 (T3) *Progressive sharpening:* In this regime, sharpness continues to increase until it reaches $\lambda^H \approx 2/\eta$
97 [15, 7]. At large effective widths or small learning rates, training may conclude before reaching this
98 threshold.

99 (T4) *Late-time dynamics (EoS):* After progressive sharpening, for MSE loss, sharpness oscillates
100 around $2/\eta$. For cross-entropy loss, the sharpness oscillates when reaching approximately $2/\eta$, while
101 decreasing over longer time scales [7].

102 In this work, we show that the sharpness dynamics heavily depends on the initialization and param-
103 eterization of the network and not every training trajectory shows all four regimes. For instance,
104 Figure 1(b, e) shows that FCNs in SP with large initialization (or large effective width) do not exhibit
105 EoS, even when loss decreases to a value below 10^{-5} . Following the early transient regime, sharpness
106 monotonically decreases, with only a nominal increase towards late training. In contrast, Figure 1(c,
107 f) shows that FCNs in μ P (or small effective width) do not experience an initial sharpness reduction at
108 small learning rates ($\eta < \eta_{\text{loss}}$). Rather, sharpness continues to increase until it reaches $2/\eta$ and then
109 oscillates around it. At large learning rates ($\eta > \eta_{\text{sharp}}$), sharpness catapults and eventually settles into
110 the same trend as above. In Appendix D, we show that these trends remain consistent when NTP is
111 used instead of SP. Given this similarity in the training dynamics between SP and NTP, we use NTP
112 for theoretical analysis for clarity and SP in realistic experiments for implementation convenience.

113 Figure 2 (and Figure 8 in Appendix E.5) demonstrates that the UV model displays all four training
114 regimes. It also captures the cases where sharpness reduction or EoS is not observed. Therefore, the
115 simplified UV model can serve as an effective model for understanding these universal behaviors in
116 the sharpness dynamics. In the subsequent section, we perform fixed point analysis of the UV model
117 and probe the origin of these complex phenomena in later sections.

118 4 Fixed Point Analysis of the UV model

119 Under GD, the parameters of the UV model are updated as $U_{t+1} = U_t - \eta \frac{\Delta f_t \mathbf{v}_t \mathbf{x}^T}{\sqrt{n_{\text{eff}}}}$, $\mathbf{v}_{t+1} =$
 120 $\mathbf{v}_t - \eta \frac{\Delta f_t U_t \mathbf{x}}{\sqrt{n_{\text{eff}}}}$, where η is the learning rate and $\Delta f_t := f(\mathbf{x}; \theta_t) - y$ is the residual at training step
 121 t . In function space, the dynamics can be completely described using the residual Δf_t and trace of
 122 the loss Hessian $\lambda := \text{Tr } H = \frac{1}{n_{\text{eff}}} (\mathbf{x}^T U^T U \mathbf{x} + \mathbf{v}^T \mathbf{v} \mathbf{x}^T \mathbf{x})$, which is also the scalar neural tangent
 123 kernel in this case. The function space dynamics of the UV model can be fully described using two
 124 coupled non-linear equations (for derivation, see Appendix E.1):

$$\Delta f_{t+1} = \Delta f_t \left(1 - \eta \lambda_t + \frac{\eta^2 \|\mathbf{x}\|^2}{n_{\text{eff}}} \Delta f_t (\Delta f_t + y) \right), \quad (1)$$

$$\lambda_{t+1} = \lambda_t + \frac{\eta \|\mathbf{x}\|^2}{n_{\text{eff}}} \Delta f_t^2 \left(\eta \lambda_t - 4 \frac{(\Delta f_t + y)}{\Delta f_t} \right), \quad (2)$$

125 with effectively three parameters: η , $\|\mathbf{x}\|/\sqrt{n_{\text{eff}}}$ and y . While similar equations have been considered
 126 in previous works [22, 36, 1], the generalization to generic parameterizations is novel and would be
 127 crucial in observing different sharpness phenomena such as EoS. The $y = 0$ case has been analyzed
 128 in prior works [22, 17] for understanding catapult dynamics. Here, λ can only decrease with time, as
 129 can be seen from Equation (2) with $\eta < \eta_{\text{max}} = 4/\lambda_0$ (training diverges if $\eta > \eta_{\text{max}}$). As a result, the
 130 model does not exhibit progressive sharpening and EoS. Below we focus on the case $y > 0$, which
 131 allows for λ to increase in time and consequently, much richer dynamics.

132 Equations (1) and (2) have four distinct fixed points/lines (referred to as I-IV) as detailed in Table 1
 133 of Appendix E.3. The fixed line I defines a zero-loss line, meaning $\ell = 0$ for all points in I; the points
 134 in I are stable for $\eta \lambda < 2$ and unstable otherwise. Fixed point II at $(-y, 0)$ corresponds to the origin
 135 in parameter space ($U, \mathbf{v} = 0$) and it is a saddle point of the dynamics for convergent learning rates
 136 η . Both I and II are also fixed points of the GD optimization, i.e., critical points of the loss. The
 137 loss Hessian at I is positive definite, while fixed point II is a saddle point in the loss landscape. The
 138 remaining two fixed points III and IV are unstable and exist only in function space, representing
 139 non-trivial parameter space dynamics that leave the function space dynamics invariant.

140 Figure 2 shows the fixed points and the vector field $\hat{G}(\Delta f, \lambda)$ determined by Equations (1) and (2),
 141 which illustrates the direction of the updates at each point. Note that the stability of the fixed line
 142 (I) does not follow from \hat{G} alone, as the magnitude G is required to determine stability. Figure 2
 143 also shows training trajectories for various parameter values. Using λ as a proxy for sharpness, we
 144 see there are regions where λ increases (colored yellow) and decreases (colored green) along the
 145 flow, which we refer to as progressive sharpening and sharpness reduction, respectively. It follows
 146 from Equation (2) that the condition $\eta \lambda \Delta f = 4(\Delta f + y)$ separates these regions. Importantly, the
 147 parameters η , $\|\mathbf{x}\|/\sqrt{n_{\text{eff}}}$ and y influence the position of the fixed points. This, in turn, affects the extent
 148 of different regions and the vector field \hat{G} , as illustrated in Figure 2. In particular, on decreasing
 149 effective width n_{eff} , or increasing learning rate η , fixed points III and IV move inward (see fixed point
 150 expressions in Table 1), which relatively enlarges the progressive sharpening region while shrinking
 151 the overall convergent region. Overall, these illustrations demonstrate how the local stability and
 152 relative position of the fixed points collectively impact the dynamics. In the subsequent section, we
 153 discuss the dynamics in detail.

154 5 Understanding Sharpness Dynamics in the UV model

155 In this section, we describe the origin of different robust phenomena in the dynamics of sharpness
 156 using the fixed point and linear stability analysis from the previous section. This explains the four
 157 training regimes observed in the UV model. We will discuss the influence of effective width and
 158 initializations, shedding light on the differences between NTP and μP . For simplicity, we assume
 159 $\|\mathbf{x}\| = 1$, while allowing n_{eff} to vary continuously. Note that we use $\lambda := \text{Tr } H$ from the previous
 160 section as a proxy for sharpness; we have verified that the top eigenvalue of the Hessian of the loss
 161 also follows λ (see Appendix E.5), although it is more difficult to analyze analytically.

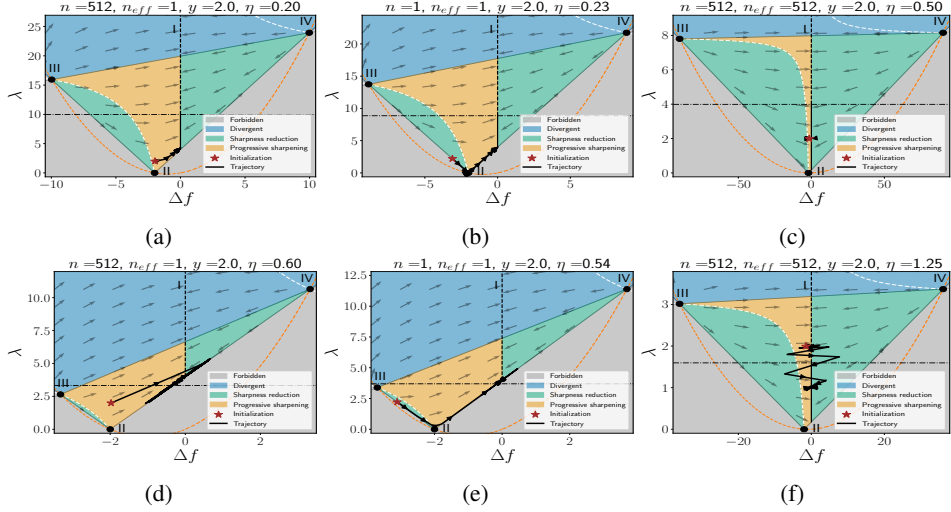


Figure 2: Training trajectories of the UV model with $\|\mathbf{x}\| = 1$ and $y = 2$ in the $(\Delta f, \lambda)$ plane for different values of n , n_{eff} and η . The columns show initializations with different n and n_{eff} , while the rows represent increasing learning rates for fixed initializations. The horizontal dash-dot line $\eta\lambda = 2$ separates the stable (solid black vertical line) and unstable (dashed black vertical line) fixed points along the zero loss fixed line I. Forbidden regions, $2\|\mathbf{x}\|\Delta f + y|\sqrt{n_{\text{eff}}}| > \lambda$, (see Appendix E.2) are shaded gray. The nullclines $\Delta f_{t+1} = \Delta f_t$ and $\lambda_{t+1} = \lambda_t$ are shown as orange and white dashed curves, respectively. Sharpness reduction, progressive sharpening, and divergent regions are colored green, yellow, and blue. The gray arrows indicate the local vector field $\hat{G}(\Delta f, \lambda)$, which is the direction of the updates. The training trajectories are depicted as black lines with arrows, with the star marking the initialization. In all cases, $\eta_c = \sqrt{n_{\text{eff}}}/2$ (introduced in Section 5.2).

162 5.1 Understanding Sharpness Trends Throughout Training

163 Figure 2 shows that the training dynamics can exhibit different behavior depending on the initial
164 region. Below we summarize these based on empirical observations.

165 (R1) *Progressive sharpening region*: As shown in Figure 2(a, d), initialization in this region experi-
166 ences an upward push due to the flow originating from fixed point II, resulting in a steady increase in
167 λ . Depending on η relative to a critical learning rate η_c (introduced in Section 5.2) different late-time
168 dynamics arises. For $\eta < \eta_c$, training converges to stable fixed points on the zero-loss line (I), as
169 shown in Figure 2(a). When $\eta > \eta_c$, all points along the zero-loss line (I) become unstable, as shown
170 in Figure 2(d). In this case, the network eventually converges to a line segment joining fixed points II
171 and IV (the EoS manifold), where it continues to oscillate indefinitely between these fixed points,
172 leading to the EoS phenomena. This will be analyzed in more depth in the subsequent section.

173 (R2) *Sharpness reduction region between fixed points II and III*: Figure 2(b, e) show that initializations
174 in this region undergo a decrease in λ as the flow is towards saddle point II. On approaching this
175 saddle point, the dynamics slows down, resulting in the intermediate saturation regime. Eventually,
176 training moves away from this saddle and enters the progressive sharpening region. From here on,
177 the dynamics becomes akin to the case (R1).

178 (R3) *Sharpness reduction region b/w fixed line I and point IV*: Initializations in this region either
179 converge to the nearby zero-loss solution for ($\eta < \eta_c$) or enter the progressive sharpening region for
180 ($\eta > \eta_c$). In the latter case, the dynamics resembles those of case (R1).

181 So far, we have described the resultant dynamics when training is initialized in different regimes.
182 Below, we describe the conditions which typically exhibit initialization these regimes.

183 **Neural Tangent Parameterization**: In NTP, Δf and λ follow normal distributions: $\Delta f_0 \sim$
184 $\mathcal{N}(-y, \sigma_w^4)$ and $\lambda_0 \sim \mathcal{N}(2\sigma_w^2, 4\sigma_w^4/n)$. Hence, the model can be initialized in any of the three
185 regions described above. Moreover, fixed points III and IV move outward with increasing width,
186 affecting the local vector field $\hat{G}(\Delta f, \lambda)$. At large widths $n \gg 1$, $\hat{G}(\Delta f_0, \lambda_0)$ at initialization points
187 along $[1 \ 0]^T$ towards the zero-loss line. For small learning rates ($\eta < 2/\lambda_0$), training exponentially
188 converges to the nearest zero-loss solution (see Figure 2(c)). Regardless of the initialization region,

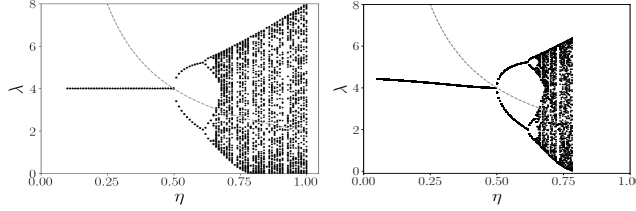


Figure 3: (left) Bifurcation diagram depicting limiting values of λ obtained by simulating Equation (3). (right) Bifurcation diagram of the UV model. In both figures, $\|\mathbf{x}\| = 1$, $y = 2$, $n_{\text{eff}} = 1$ and $\eta_c = 0.5$.

189 the change in λ is minimal, receiving $\mathcal{O}(1/n)$ updates as per Equation (2). For large learning rates
 190 ($\eta > 2/\lambda_0$), the nearby zero-loss solution becomes unstable. Consequently, training catapults to a
 191 region with smaller λ , while bouncing between fixed points III and IV. This is the catapult effect
 192 studied in [22] and Figure 2(f) demonstrates such a trajectory. By comparison, at small widths, the
 193 dynamics follows cases (R1-R3) discussed above.

194 **Maximal Update (μ P) Parameterizations:** In contrast to NTP, the position of fixed points III-
 195 IV do not change with width n , and Δf_0 follows the distribution: $\Delta f_0 \sim \mathcal{N}(-y, \sigma_w^4/n)$, while
 196 λ_0 distribution remains unchanged. Consequently, at large widths, the model is initialized at
 197 $(-y, 2\sigma_w^2)$, right above fixed point II in the progressive sharpening region (R1), satisfying the
 198 condition $\eta\lambda_0\Delta f_0 < 4(\Delta f_0 + y)$. Figure 2(a, d) shows such a trajectory. At small widths, fluctua-
 199 tions increase, making it plausible for μ P networks to start in the sharpness reduction regions. In this
 200 case, the dynamics follow case (R2) or (R3).

201 5.2 Understanding Edge of Stability

202 This section analyzes the EoS behavior in the UV model, particularly from the fixed point perspective.
 203 As discussed in the previous section, the EoS behavior in the UV model arises when all fixed points
 204 along the zero-loss line (I) become unstable wrt the learning rate. Yet, the gradients updates (shown
 205 as gray arrows in Figure 2) continue to point towards the zero loss line. As a result training is trapped
 206 in this region, converging to the line segment that joins fixed points II and IV —referred to as the
 207 EoS manifold —where it oscillates indefinitely.

208 **EoS Manifold is an Attractor:** By examining the two-step dynamics akin to Ref. [1, 5], we show
 209 in Appendix E.7 that training converges to the EoS manifold above a critical learning rate η_c . For
 210 $\eta < \eta_c$, training converges to the stable fixed points on the zero-loss line. By comparison, for $\eta > \eta_c$,
 211 all points along the zero-loss line become unstable and the EoS manifold becomes a dynamical
 212 attractor. The critical η_c for which all points on the zero-loss line become unstable thus gives a
 213 necessary condition for EoS:

214 **Result 1.** A necessary condition for the UV model to exhibit EoS is $\eta > \eta_c = \sqrt{n_{\text{eff}}}/\|\mathbf{x}\|_y$ (see
 215 Appendix E.8 for details). It is useful to scale the learning rate as $\eta = c/\lambda_0$, in which case this
 216 condition becomes $\lambda_0 < c\|\mathbf{x}\|_y/\sqrt{n_{\text{eff}}}$. For learning rates $\eta > 2/\lambda_0$, training can catapult to regions
 217 with $\lambda_T < \lambda_0$. In such cases, the condition $\lambda_T < c\|\mathbf{x}\|_y/\sqrt{n_{\text{eff}}}$ also applies.

218 **Dynamics on the EoS Manifold and Route to Chaos:** The dynamics on the EoS manifold satisfies
 219 $\lambda = 2\|\mathbf{x}\|(\Delta f + y)/\sqrt{n_{\text{eff}}}$, coupling Δf and λ together. This yields the map $\Delta f_{t+1} = M_f(\Delta f_t)$
 220 describing the dynamics on the EoS manifold, with M_f defined as

$$M_f(\Delta f_t) := \Delta f_t + \frac{\eta\Delta f_t}{\eta_c y} \left(\frac{\eta\Delta f_t}{\eta_c y} - 2 \right) (\Delta f_t + y). \quad (3)$$

221 Figure 3(left) shows the limiting values of λ (i.e. the values of λ that the network jumps between at
 222 late times) as a function of learning rate, obtained by simulating Equation (3). We refer to this as the
 223 bifurcation diagram. As mentioned before, for $\eta > \eta_c$, the zero-loss solution becomes unstable with
 224 λ oscillating around $2/\eta$ instead of converging. These fluctuations exhibit a fractal structure, as the
 225 system undergoes a series of period-doubling transitions with an increasing learning rate. This is the
 226 well-known *period-doubling route to chaos* [25]. Figure 3(right) shows the bifurcation diagram of
 227 the UV model for $y = 2$. The bifurcation diagram extends up to $\eta \approx 0.8$ before diverging at
 228 higher learning rates. This leads us to the following corollary of Result 1.

229 **Corollary 5.1.** Let η_{max} be the maximum trainable learning rate for a given initialization. The
 230 bifurcation diagram is observed up to $\eta < \eta_{\text{max}}$. If $\eta_{\text{max}} < \eta_c$, the UV model does not exhibit EoS.

231 These results suggest that models with small λ_0 and n_{eff} are more prone to show EoS behavior. As a
232 result, μP networks or those with small initial weight variance are more likely to exhibit EoS. On
233 the other hand, large-width NTP networks may not show EoS behavior at all. In Appendix A, we
234 validate this prediction in real-world scenarios.

235 **Connections to sub-quadratic loss:** Ref. [23] demonstrated that GD on sub-quadratic loss with
236 large learning rates inherently results in EoS behavior. Here, we show that the loss on the EoS
237 manifold of the UV model is sub-quadratic near its minimum. As noted above, the dynamics on
238 the EoS manifold satisfies $\lambda = 2\|\mathbf{x}\|(\Delta f + y)/\sqrt{n_{\text{eff}}}$. The loss on the EoS manifold is then given by
239 $\mathcal{L}(\theta) = \frac{1}{2}\Delta f^2 = \frac{y^2}{2}(\frac{\eta\lambda}{2} - 1)^2$, where θ denotes the parameters. Since $\lambda \sim \mathcal{O}(\|\theta\|^2)$, the loss is
240 of the form $\mathcal{L}(\theta) \approx \frac{1}{2}(a\|\theta\|^2 - b)^2$ and is sub-quadratic near its minimum. The GD dynamics near
241 the minimum is given by a cubic map, which is known to show the period-doubling route to chaos
242 [26]. Ref. [6] showed a similar route to chaos by considering a two-layer network with quadratic
243 activation, with the last layer vector \mathbf{v} fixed through training and each entry set to one. In this model,
244 the loss is sub-quadratic by *construction* ($(\|U\mathbf{x}\|^2 - y)^2$) and the dynamics is given by a cubic map.

245 6 Discussion

246 The applicability of the fixed point analysis extends well beyond the UV model and can be employed
247 in settings involving complex architectures and adaptive optimizers. A prerequisite for applying this
248 method is the closure of the dynamical equations describing the model. By analyzing the fixed points
249 of such equations in broader classes of models, we can gain significant insights into their training
250 dynamics, thereby advancing our understanding of non-convex optimization in neural networks.

251 Various results such as the phase diagram of EoS, the bifurcation diagram, and the late-time sharpness
252 analysis depend on the training time. Nevertheless, we found that training the models longer does
253 not impact the conclusions presented. In Appendix F, we show our results are robust for reasonably
254 small batch sizes ($B \approx 512$). For even smaller batch sizes, the dynamics becomes noise-dominated,
255 and separating the inherent dynamics from noise becomes challenging.

256 References

- 257 [1] Atish Agarwala, Fabian Pedregosa, and Jeffrey Pennington. A second-order regression model
258 shows edge of stability behavior. In *OPT 2022: Optimization for Machine Learning (NeurIPS
259 2022 Workshop)*, 2022.
- 260 [2] Kwangjun Ahn, Sébastien Bubeck, Sinho Chewi, Yin Tat Lee, Felipe Suarez, and Yi Zhang.
261 Learning threshold neurons via the "edge of stability". *ArXiv*, abs/2212.07469, 2022.
- 262 [3] Sanjeev Arora, Zhiyuan Li, and Abhishek Panigrahi. Understanding gradient descent on
263 the edge of stability in deep learning. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song,
264 Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International
265 Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*,
266 pages 948–1024. PMLR, 17–23 Jul 2022.
- 267 [4] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal
268 Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao
269 Zhang. JAX: composable transformations of Python+NumPy programs, 2018.
- 270 [5] Lei Chen and Joan Bruna. Beyond the edge of stability via two-step gradient updates. *arXiv*,
271 2206.04172, 2023.
- 272 [6] Xuxing Chen, Krishnakumar Balasubramanian, Promit Ghosal, and Bhavya Agrawalla. From
273 stability to chaos: Analyzing gradient descent dynamics in quadratic regression. 2310.01687,
274 2023.
- 275 [7] Jeremy Cohen, Simran Kaur, Yuanzhi Li, J Zico Kolter, and Ameet Talwalkar. Gradient descent
276 on neural networks typically occurs at the edge of stability. In *International Conference on
277 Learning Representations*, 2021.

- 278 [8] Alex Damian, Eshaan Nichani, and Jason D. Lee. Self-stabilization: The implicit bias of
279 gradient descent at the edge of stability. In *The Eleventh International Conference on Learning*
280 *Representations*, 2023.
- 281 [9] Li Deng. The mnist database of handwritten digit images for machine learning research. *IEEE*
282 *Signal Processing Magazine*, 29(6):141–142, 2012.
- 283 [10] Darshil Doshi, Tianyu He, and Andrey Gromov. Critical initialization of wide and deep neural
284 networks through partial jacobians: General theory and applications. *arXiv*, 2111.12143, 2023.
- 285 [11] Gintare Karolina Dziugaite and Daniel M Roy. Computing nonvacuous generalization bounds
286 for deep (stochastic) neural networks with many more parameters than training data. *arXiv*,
287 1703.11008, 2017.
- 288 [12] Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
289 recognition. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*,
290 pages 770–778, 2016.
- 291 [13] Jonathan Heek, Anselm Levskaya, Avital Oliver, Marvin Ritter, Bertrand Rondepierre, Andreas
292 Steiner, and Marc van Zee. Flax: A neural network library and ecosystem for JAX, 2020.
- 293 [14] Arthur Jacot, Franck Gabriel, and Clement Hongler. Neural tangent kernel: Convergence and
294 generalization in neural networks. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman,
295 N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*,
296 volume 31. Curran Associates, Inc., 2018.
- 297 [15] Stanislaw Jastrzebski, Maciej Szymczak, Stanislav Fort, Devansh Arpit, Jacek Tabor,
298 Kyunghyun Cho*, and Krzysztof Geras*. The break-even point on optimization trajectories of
299 deep neural networks. In *International Conference on Learning Representations*, 2020.
- 300 [16] Yiding Jiang, Behnam Neyshabur, Hossein Mobahi, Dilip Krishnan, and Samy Bengio. Fantastic
301 generalization measures and where to find them. *arXiv*, 1912.02178, 2019.
- 302 [17] Dayal Singh Kalra and Maissam Barkeshli. Phase diagram of early training dynamics in deep
303 neural networks: effect of the learning rate, depth, and width. In *Thirty-seventh Conference on*
304 *Neural Information Processing Systems*, 2023.
- 305 [18] Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping
306 Tak Peter Tang. On large-batch training for deep learning: Generalization gap and sharp minima.
307 *ArXiv*, 1609.04836, 2016.
- 308 [19] Lingkai Kong and Molei Tao. Stochasticity of deterministic gradient descent: Large learning
309 rate for multiscale objective function. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan,
310 and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages
311 2625–2638. Curran Associates, Inc., 2020.
- 312 [20] Itai Kreisler, Mor Shpigel Nacson, Daniel Soudry, and Yair Carmon. Gradient descent monotonically
313 decreases the sharpness of gradient flow solutions in scalar networks and beyond. *ArXiv*,
314 2305.13064, 2023.
- 315 [21] Alex Krizhevsky. Learning multiple layers of features from tiny images. 2009.
- 316 [22] Aitor Lewkowycz, Yasaman Bahri, Ethan Dyer, Jascha Sohl-Dickstein, and Guy Gur-Ari. The
317 large learning rate phase of deep learning: the catapult mechanism. *ArXiv*, 2003.02218, 2020.
- 318 [23] ChaoKunin Ma, Lie Wu, and Lexing Ying. Beyond the quadratic approximation: The multiscale
319 structure of neural network loss landscapes. *Journal of Machine Learning*, 1(3):247–267, 2022.
- 320 [24] Lorenzo Noci, Alexandru Meterez, Thomas Hofmann, and Antonio Orvieto. Why do learning
321 rates transfer? reconciling optimization and scaling limits for deep learning, 2024.
- 322 [25] Edward Ott. *Chaos in Dynamical Systems*. Cambridge University Press, 2 edition, 2002.
- 323 [26] Thomas D. Rogers and David C. Whitley. Chaos in the cubic mapping. *Mathematical Modelling*,
324 4(1):9–25, 1983.

- 325 [27] Mihaela Rosca, Yan Wu, Chongli Qin, and Benoit Dherin. On a continuous time model of
326 gradient descent dynamics and instability in deep learning. *Transactions on Machine Learning*
327 *Research*, 2023.
- 328 [28] Vaishaal Shankar, Alexander W. Fang, Wenshuo Guo, Sara Fridovich-Keil, Ludwig Schmidt,
329 Jonathan Ragan-Kelley, and Benjamin Recht. Neural kernels without tangents. In *ICML*, 2020.
- 330 [29] Jascha Sohl-Dickstein, Roman Novak, Samuel S. Schoenholz, and Jaehoon Lee. On the infinite
331 width limit of neural networks with a standard parameterization. *arXiv*, 2001.07301, 2020.
- 332 [30] Minhak Song and Chulhee Yun. Trajectory alignment: Understanding the edge of stability
333 phenomenon via bifurcation theory. *arXiv*, 2307.04204, 2023.
- 334 [31] Yuqing Wang, Zhenghao Xu, Tuo Zhao, and Molei Tao. Good regularity creates large learning
335 rate implicit biases: edge of stability, balancing, and catapult. *arXiv*, 2310.17087, 2023.
- 336 [32] Zixuan Wang, Zhouzi Li, and Jian Li. Analyzing sharpness along GD trajectory: Progressive
337 sharpening and edge of stability. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and
338 Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022.
- 339 [33] Jingfeng Wu, Vladimir Braverman, and Jason D. Lee. Implicit bias of gradient descent for
340 logistic regression at the edge of stability. *arXiv*, 2305.11788, 2023.
- 341 [34] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for
342 benchmarking machine learning algorithms, 2017. cite arxiv:1708.07747Comment: Dataset is
343 freely available at <https://github.com/zalando-research/fashion-mnist> Benchmark is available at
344 <http://fashion-mnist.s3-website.eu-central-1.amazonaws.com/>.
- 345 [35] Greg Yang and Edward J. Hu. Tensor programs iv: Feature learning in infinite-width neural
346 networks. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International*
347 *Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*,
348 pages 11727–11737. PMLR, 18–24 Jul 2021.
- 349 [36] Libin Zhu, Chaoyue Liu, Adityanarayanan Radhakrishnan, and Mikhail Belkin. Quadratic
350 models for understanding neural network dynamics. *ArXiv*, 2205.11787, 2022.
- 351 [37] Xingyu Zhu, Zixuan Wang, Xiang Wang, Mo Zhou, and Rong Ge. Understanding edge-of-
352 stability training dynamics with a minimalist example. In *The Eleventh International Conference*
353 *on Learning Representations*, 2023.

354 A Predictions and Verifications in Real-world Scenarios

355 The preceding analysis offers broader insights and predictions for optimization in real-world models.
356 In this section, we study realistic architectures with real and synthetic datasets and examine the extent
357 to which insights from the UV model generalize.

358 **Experimental Setup:** Consider a network $f(\mathbf{x}; \theta)$, with trainable parameters θ , initialized using
359 normal distribution with zero mean and variance σ_w^2 in appropriate parameterization. In this section,
360 we use the interpolating parameterization with $s \in [0, 1]$ (detailed in Appendix C.2.1), where networks
361 with $s = 0$ are equivalent to networks in SP as width n goes to infinity and those with $s = 1$ are in μ P.
362 The network is trained on a dataset \mathcal{D} with P examples using MSE loss and GD. The learning rate is
363 scaled as $\eta = c/\lambda_0^H$, where c is the learning rate constant, and λ_0^H is the sharpness at initialization.
364 Additional details provided in figure captions and Appendix C.2.

365 **Implications of Initialization and Parameterization for Real-world Models:** The analysis in
366 Section 5.1 unveils crucial insights into the implicit biases of parameterization in real-world networks.
367 Figure 2(a, d) shows that μ P networks begin training in a flat region of the landscape, where
368 gradients point towards increasing sharpness, and approach the zero loss line while maintaining
369 a low sharpness bias. In contrast, networks in NTP (or equivalently SP), characterized as large
370 initializations, experience sharpness reduction during early training and might not approach with a
371 minimal sharpness bias. The agreement of these observations with networks trained on real-world
372 datasets (Figure 1) suggests that these inherent biases hold in practical scenarios.

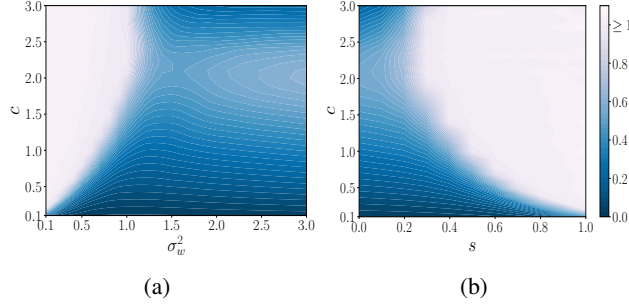


Figure 4: (a) Heatmap of $\eta\bar{\lambda}^H/2$ of ReLU FCNs with $s = 0$ trained on a 5k subset of CIFAR-10 for 10k steps, with the weight variance σ_w^2 and learning rate multiplier $c = \eta\lambda_0^H$ as axes. $\bar{\lambda}^H$ is obtained by averaging λ_t^H over last 200 steps. As the color varies from blue to white, $\eta\bar{\lambda}^H/2$ increases, where the brightest white region indicates the EoS regime with $\eta\bar{\lambda}^H/2 \geq 1$. (b) Same heatmap with fixed

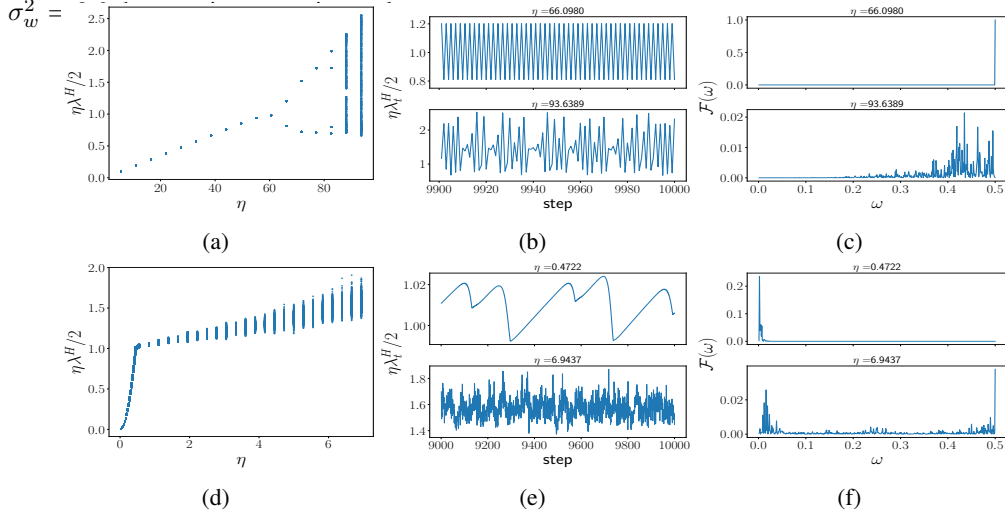


Figure 5: 2-layer linear FCNs trained on (first row) 5,000 iid random examples with unit output dimension and (second row) 5,000 CIFAR-10 examples. Different columns correspond to the bifurcation diagram, late-time sharpness trajectories, and the power spectrum of sharpness trajectories. The power spectrum is computed using the last 1000 steps of the trajectories.

373 **Sharpness & Weight-Norm Correlation and the Origin of Four Regimes:** Section 5.1 revealed
 374 that, for a wide variety of initializations, at early times trajectories move closer to the saddle point
 375 II, resulting in an interim decrease in λ (also proportional to the weight norm in this case), before
 376 eventually increasing. This critical point where all parameters are zero also exists in real-world
 377 models. We thus anticipate that in real-world models, the origin of the four training regimes may
 378 be related to a similar mechanism. This would predict a decrease in weight norm as training passes
 379 near the saddle point, followed by an eventual increase. In Appendix G, we validate this hypothesis.
 380 During the sharpness reduction and intermediate saturation regimes, we see a decrease in the weight
 381 norm, followed by an increase in the weight norm as the network undergoes progressive sharpening,
 382 following the prediction from the UV model.

383 **The Phase Diagram of Edge of Stability:** Result 1 presents a necessary condition for EoS to
 384 occur in the UV model: $\lambda_0 < c\|x\|y/\sqrt{n_{\text{eff}}}$. In real-world models, the initial sharpness λ_0^H can be
 385 controlled using the initial variance of the weights σ_w^2 . Therefore, this result predicts that real-world
 386 models with (i) small initial weight variance σ_w^2 , (ii) large interpolating parameter s , or (iii) large
 387 learning constant c are more likely to exhibit EoS behavior. Figure 4 shows the phase diagram of
 388 EoS, validating these predictions. Additional phase diagrams in Appendix H indicate an enhanced
 389 tendency for CNNs and ResNets to exhibit EoS.

390 **Route to Chaos and Bifurcation Diagrams:** The analysis in Section 5.2 unveiled structured
 391 fluctuations in λ at the EoS, with a period-doubling route to chaos observed as the learning rate is

392 tuned. This motivates us to analyze fluctuations at the EoS in real-world models trained on realistic
 393 and synthetic datasets. Figure 5 shows the bifurcation diagram, late-time sharpness trajectories, and
 394 power spectrum of sharpness trajectories for a 2-layer linear FCN. In the first row, the model is
 395 trained on random synthetic data with 5,000 iid examples with unit output dimension, whereas, in
 396 the second row, on a 5,000 example subset of CIFAR-10. Similar to the UV model, FCNs trained
 397 on random data exhibit a period-doubling route to chaos, as shown in Figure 5(a). By comparison,
 398 FCNs trained on CIFAR-10 only show dense bands in the sharpness rather than exhibiting a clear
 399 period-doubling route to chaos.

400 On analyzing the sharpness trajectories at EoS, we observe long-range correlations in time in real
 401 datasets, with fluctuations increasing with the learning rate (see Figure 5(e)). By comparison,
 402 sharpness trajectories of models trained on random datasets exhibit short-period oscillations (see
 403 Figure 5(b)). The power spectrum of these sharpness trajectories further quantifies these observations,
 404 as shown in Figure 5(c, f). In the random dataset case, high-frequency modes corresponding to the
 405 period-doubling route to chaos emerge at EoS as shown in Figure 5(c). In contrast, real datasets
 406 exhibit low-frequency modes at small learning rates. As the learning rate is increased, high-frequency
 407 modes, reminiscent of the period-doubling route to chaos, start emerging (see Figure 5(f)). In
 408 Appendix I.1, we demonstrate that CNNs and ResNets trained on image datasets show dense bands
 409 of sharpness similar to those in FCNs.

410 To understand when the period-doubling route to chaos arises, we perform further analysis in
 411 Appendix I. A key determining feature appears to be whether the singular value spectrum of the
 412 input-input and output-input covariance matrices are flat or have power-law decay. In Appendix I.2,
 413 we show that a 2-layer FCNs trained on a random dataset with power-law singular value spectrum
 414 in the input exhibits dense sharpness bands. In Appendix I.3 we show that linear FCNs trained on
 415 synthetic datasets with random inputs, such as teacher-student settings and generative settings (details
 416 in Appendix C.1), exhibit the period-doubling route to chaos. In contrast, non-linear networks trained
 417 on these tasks exhibit dense sharpness bands as observed in real datasets. These observations shed
 418 some light on the nature of EoS observed in realistic settings. Nevertheless, a complete understanding
 419 of sharpness fluctuations at EoS requires a separate detailed examination.

420 B Further Discussion on Related Works

421 [22] examined the early training dynamics of wide networks at large learning rates. Using the top
 422 eigenvalue of the Neural Tangent Kernel (NTK) λ^K at initialization ($t = 0$), they revealed a ‘catapult
 423 phase’, $2/\lambda_0^K < \eta < \eta_{\max}$, in which training converges despite an initial spike in training loss. [17]
 424 analyzed early training dynamics for arbitrary depths and width and revealed a ‘sharpness reduction
 425 phase’, $2/\lambda_0^H < \eta < c_{\text{loss}}/\lambda_0^H$, which opens up significantly as c_{loss} increases with depth and $1/\text{width}$.

426 Beyond early training, sharpness continues to increase, until it reaches a break-even point [15], beyond
 427 which GD dynamics typically enters the EoS regime [7]. This has motivated various theoretical
 428 studies to understand GD dynamics at large learning rates: [3, 27, 37, 33, 5, 2, 20, 30, 6]. In particular,
 429 Ma et al. [23] showed that loss functions with sub-quadratic growth exhibit EoS behavior. [3] show
 430 that normalized gradient descent reaches the EoS regime. [8] analyze the dynamics of the cubic
 431 approximation of the loss. Assuming a negative correlation between the gradient direction and the
 432 top eigenvector of Hessian, they show that gradient descent dynamics enters a stable cycle in the EoS
 433 regime. [2] analyze EoS in a single-neuron 2-layer network and a simplified three-parameter ReLU
 434 network assuming the existence of a ‘forward invariant subset’ near the minima. [20] analyzed scalar
 435 linear networks to show that the sharpness attained by the gradient flow dynamics monotonically
 436 decreases in the EoS regime. [33] demonstrate that gradient descent, with any learning rate in the
 437 EoS regime, optimizes logistic regression with linearly separated data over large time scales. Below,
 438 we discuss closely related works in detail and clarify the distinction with our work.

439 [32] analyze EoS in a 2-layer linear network using the norm of the last layer. They solely focus on
 440 cases that exhibit progressive sharpening right from initializations by considering assumptions (refer
 441 to Assumptions 4.1 and 4.2 of their paper) on the training dataset. Contrary to these assumptions, Fig-
 442 ure 1(e) demonstrates that such assumptions are invalid in many realistic settings, where progressive
 443 sharpening is not observed at all.

444 [1] showed that a modified model exhibits progressive sharpening and two-step oscillations at
 445 EoS using NTK as the proxy. They state that the UV model does not exhibit EoS behavior (see

446 Section 3.2.1 of the referenced paper). This is because their analysis is restricted to the Standard
 447 Parameterization corresponding to Figure 2(c, f) (c.f. Figure 8 of [1]). In contrast, we show that
 448 the UV model exhibits EoS behavior under the appropriate choice of parameterization and training
 449 example.

450 [37] proved EoS convergence for the loss $\frac{1}{4}(x^2y^2 - 1)^2$, where $x, y \in \mathbb{R}$. Additionally, they
 451 empirically demonstrated a bifurcation diagram in the space of abstract variables of x and y . It is
 452 worth noting that while these bifurcations arise from the same underlying behavior, they contrast with
 453 our route to chaos bifurcation diagrams which quantify sharpness fluctuations with learning rate.

454 [5] analyze two-step gradient updates of a single-neuron network and matrix factorization to gain
 455 insights into EoS. Similar to our work, they show a bifurcation diagram of sharpness against the
 456 learning rate for the matrix factorization problem. While the scalar matrix factorization problem can
 457 be mapped to the UV model with a specific choice of $\frac{\|x\|}{\sqrt{n_{\text{eff}}}}$, it is not straightforward to apply their
 458 conclusions to the neural network setting, as it requires the correct choice of parameterization. In
 459 particular, the UV model under NTP parameterization, as shown in Figure 2(c, f), does not display
 460 EoS behavior at considerable widths, a finding also noted by [1]. Observing EoS requires the correct
 461 choice of the parameterization (μP) and training example. Furthermore, although the scalar matrix
 462 factorization in [5] can be mapped to a special case of the UV model considered in our work, we
 463 provide significant additional insights. In particular, with respect to the bifurcation phenomena in the
 464 UV model, we explain the existence of an attractor submanifold on which the EoS behavior occurs.
 465 We further show that on the EoS submanifold, the loss becomes subquadratic in nature and the
 466 gradient descent dynamics therefore become approximated by the cubic map, which is well-studied in
 467 the chaos literature. This makes clear the origin of the period-doubling route to chaos. Additionally,
 468 in Section 6, we extend the analysis of EoS beyond the UV model, comparing sharpness trajectories
 469 of synthetic and real datasets at EoS. In contrast to the synthetic datasets, sharpness trajectories
 470 of real datasets show long-range correlations in time. We take the first steps by attributing these
 471 long-range correlations to correlations in the dataset. Note that this setting cannot be mapped to the
 472 matrix factorization setting.

473 [30] show that late-time trajectories oscillate around $2f/\eta\ell'$, where f is the network output and ℓ' is
 474 the derivative of the loss. They refer to the term bifurcation diagram to describe these phenomena,
 475 contrasting with the sharpness versus learning rate bifurcation diagrams presented in our study. We
 476 quote Section 3 from their paper ".we plot the bifurcation diagram $q = r(p) = \ell'(p)/p$ and observe
 477 that GD trajectories tend to align with this curve." Here, p and q correspond to Δf and $\frac{2}{\eta\lambda}$ in our
 478 setting. They plot trajectories in the $(p, q) \equiv (\Delta f, \frac{2}{\eta\lambda})$ plane and their condition $q = r(p)$ simply
 479 corresponds to the EoS condition $\lambda = \frac{2}{\eta}$ for MSE loss. In contrast, our work presents bifurcation
 480 diagrams resulting from how the sharpness fluctuations vary with the learning rate. Therefore, the
 481 bifurcation diagrams from these works are not directly related to the route-to-chaos bifurcation
 482 diagrams presented in our work.

483 [19] were the first ones to show that gradient descent dynamics becomes chaotic at large learning
 484 rates and converges to a statistical distribution instead of a minimum.

485 [6] analyzed large learning rate dynamics of toy models which are characterized by a one-dimensional
 486 cubic map and demonstrated five different training phases: (a) monotonic, (b) catapult, (c) periodic,
 487 (d) chaotic, and (e) divergent. In particular, they considered a two-layer network with quadratic
 488 activation, where the last layer vector v is not trained and each entry is set to one. This model belongs
 489 to a family that is effectively described by one variable Δf . In this model, the loss is sub-quadratic
 490 by construction $(\|Ux\|^2 - y)^2$. In contrast, the UV model that we study is an effectively two-variable
 491 model and in these cases, training dynamically finds the attractive EoS manifold such that the loss
 492 has a sub-quadratic nature on this submanifold.

493 Concurrent work by [31] categorizes training trajectories into three stages: (i) sharpness reduction,
 494 (ii) progressive sharpening, and (iii) edge of stability. They argue that different large learning rate
 495 behavior depends on the 'regularity' of the loss landscape. Specifically, they generalize toy landscapes
 496 from existing studies with parameters controlling the regularity. They show that models with good
 497 regularity first experience a decrease in sharpness and then progressive sharpening and enter the edge
 498 of stability.

499 [24] examined the sharpness dynamics of networks with with parameterization. They argue that
 500 the learning rate transfer property of $\mu\mathbf{P}$ is correlated with consistent sharpness trajectories
 501 across varying depths and widths.

502 C Experimental details

503 C.1 Datasets

504 **Standard image datasets:** We considered the MNIST [9], Fashion-MNIST [34], and CIFAR-10
 505 [21] datasets. The images are standardized to have zero mean and unit variance across the feature
 506 dimensions, and target labels are represented as one-hot encodings.

507 **Random dataset:** We construct a random dataset $(X, Y) = \{(\mathbf{x}^\mu, \mathbf{y}^\mu)\}_{\mu=1}^P$ with $\mathbf{x}^\mu \sim \mathcal{N}(0, I)$
 508 and $\mathbf{y}^\mu \sim \mathcal{N}(0, I)$, both sampled independently. Note there is no correlation between inputs and
 509 outputs.

510 **Teacher-student dataset:** Consider a teacher network $f(\mathbf{x}; \theta_0)$ with θ_0 initialized randomly as
 511 described in Appendix C.2. Then, we construct a teacher-student dataset $(X, Y) = \{(\mathbf{x}^\mu, \mathbf{y}^\mu)\}_{\mu=1}^P$
 512 with $\mathbf{x}^\mu \sim \mathcal{N}(0, I)$ and $\mathbf{y}^\mu = f(\mathbf{x}^\mu; \theta_0)$.

513 **Random power-law dataset:** Starting with the random dataset (X', Y') , we utilize the singular
 514 value decomposition of the input and output matrices

$$X' = P_x S_{x'} Q_x^T, \quad Y' = P_y S_{y'} Q_y^T. \quad (4)$$

515 Next, we rescale the k^{th} singular value of $S_{x'}$ and $S_{y'}$ as

$$(S_x)_k = A_x (S_{x'})_k k^{-B_x} \quad (S_y)_k = A_y (S_{y'})_k k^{-B_y}, \quad (5)$$

516 and re-construct input and output matrices as below

$$X = P_x S_x Q_x^T, \quad Y = P_y S_y Q_y^T. \quad (6)$$

517 The variables A_x, B_x, A_y , and B_y uniquely characterize the dataset.

518 **Generative image dataset:** Given a pre-trained network $f(\mathbf{x}; \theta)$ on a standard image dataset listed
 519 above, we construct a generative image dataset $(X, Y) = \{(\mathbf{x}^\mu, \mathbf{y}^\mu)\}_{\mu=1}^P$ with $\mathbf{x}^\mu \sim \mathcal{N}(0, I)$ and
 520 $\mathbf{y}^\mu = f(\mathbf{x}^\mu; \theta)$.

521 C.2 Models

522 **FCNs:** We considered ReLU FCNs without bias with uniform hidden layer width n .

523 **CNNs:** We considered Myrtle family ReLU CNNs [28] without any bias with a fixed number of
 524 channels in each layer, which we refer to as the width of the network.

525 **ResNets:** We adapted ResNet [12] implementations from Flax examples. Our implementation
 526 uses Layer norm and initialize the weights as $\mathcal{N}(0, \sigma_w^2/\text{fan}_{in})$. For ResNets, we refer to the number of
 527 channels in the first block as the width.

528 We implemented all models using the JAX [4], and Flax libraries [13].

529 C.2.1 Details of network parameterization

530 In this section, we describe different parameterizations used in the paper. For simplicity, we describe
 531 the parameterizations for FCNs. Nevertheless, these arguments generalize to other architectures.

532 **Standard Parameterization (SP):** Consider a neural network $f : \mathbb{R}^{d_{\text{in}}} \rightarrow \mathbb{R}^{d_{\text{out}}}$ with d layers and
 533 constant width n . Then, standard parameterization is defined as follows:

$$\begin{aligned} \mathbf{h}^{(1)}(\mathbf{x}) &= W^{(1)}\mathbf{x}, \\ \mathbf{h}^{(l+1)}(\mathbf{x}) &= W^{(l+1)}\phi\left(\mathbf{h}^{(l)}(\mathbf{x})\right), \\ \mathbf{f}(\mathbf{x}; \theta) &= W^{(d)}\phi\left(\mathbf{h}^{(d-1)}(\mathbf{x})\right), \end{aligned} \quad (7)$$

534 where $W^{(1)} \sim \mathcal{N}(0, \sigma_w^2/d_{\text{in}})$, $W^{(l)} \sim \mathcal{N}(0, \sigma_w^2/n)$ for $1 < l < d$, and $W^{(d)} \sim \mathcal{N}(0, 1/n)$; $\phi(\cdot)$ is the
 535 elementwise activation function. The input is normalized such that $\|\mathbf{x}\|^2 = d_{\text{in}}$.

536 **Neural Tangent Parameterization (NTP):** Consider a neural network $f : \mathbb{R}^{d_{\text{in}}} \rightarrow \mathbb{R}^{d_{\text{out}}}$ with d
 537 layers and constant width n . Then, the Neural Tangent Parameterization is defined as follows:

$$\begin{aligned} \mathbf{h}^{(1)}(\mathbf{x}) &= \frac{\sigma_w}{\sqrt{d_{\text{in}}}} W^{(1)}\mathbf{x}, \\ \mathbf{h}^{(l+1)}(\mathbf{x}) &= \frac{\sigma_w}{\sqrt{n}} W^{(l+1)}\phi\left(\mathbf{h}^{(l)}(\mathbf{x})\right), \\ \mathbf{f}(\mathbf{x}; \theta) &= \frac{1}{\sqrt{n}} W^{(d)}\phi\left(\mathbf{h}^{(d-1)}(\mathbf{x})\right), \end{aligned} \quad (8)$$

538 where $W^{(l)} \sim \mathcal{N}(0, 1)$ for $1 \leq l \leq d$ and $\phi(\cdot)$ is the elementwise activation function. The input is
 539 normalized such that $\|\mathbf{x}\|^2 = d_{\text{in}}$.

540 Both SP and NTP are closely related parameterizations, as constant width networks in SP with
 541 learning rate $\eta = \Theta(1/n)$ learning rate are equivalent to those in NTP [35].

542 **Interpolating Parameterization:** Consider a neural network $f : \mathbb{R}^{d_{\text{in}}} \rightarrow \mathbb{R}^{d_{\text{out}}}$ with d layers and
 543 constant width n . Let $W^{(l)}$ denote the weight matrix at layer l . Then, ‘‘interpolating parameterization’’
 544 is defined as follows:

$$\begin{aligned} \mathbf{h}^{(1)}(\mathbf{x}) &= n^{s/2} W^{(1)}\mathbf{x}, \\ \mathbf{h}^{(l+1)}(\mathbf{x}) &= W^{(l+1)}\phi\left(\mathbf{h}^{(l)}(\mathbf{x})\right), \\ \mathbf{f}(\mathbf{x}; \theta) &= \frac{1}{n^{s/2}} W^{(d)}\phi\left(\mathbf{h}^{(d-1)}(\mathbf{x})\right), \end{aligned} \quad (9)$$

545 Here, s is a parameter that interpolates between standard-like parameterization and maximal update pa-
 546 rameterization. The weight matrices are sampled from Gaussian distributions: $W^{(1)} \sim \mathcal{N}(0, \sigma_w^2/n^s)$,
 547 $W^{(l)} \sim \mathcal{N}(0, \sigma_w^2/n)$ for $1 < l < d$, and $W^{(d)} \sim \mathcal{N}(0, 1/n)$. We normalize the input such that
 548 $\|\mathbf{x}\| = 1$.

549 **Maximal update Parameterization (μP):** The maximal update parameterization corresponds to
 550 the $s = 1$ case in the above setting.

551 C.3 Details of Figures

552 **Figure 1:** Training loss and sharpness trajectories of 4-layer ReLU FCNs with $n = 512$, trained on
 553 a subset of 5,000 CIFAR-10 examples using MSE loss and GD: (a, d) SP with $\sigma_w^2 = 0.5$, (b, e) SP
 554 with $\sigma_w^2 = 2.0$, (c, f) μP with $\sigma_w^2 = 2.0$.

555 **Figure 2** Training trajectories of the UV model with $\|\mathbf{x}\| = 1$ and $y = 2$ in the $(\Delta f, \lambda)$ plane for
 556 different values of n , n_{eff} and η . The columns show initializations with different n and n_{eff} , while the

557 rows represent increasing learning rates for fixed initializations. The horizontal dash-dot line $\eta\lambda = 2$
558 separates the stable (solid black vertical line) and unstable (dashed black vertical line) fixed points
559 along the zero loss fixed line I. Forbidden regions, $2\|\mathbf{x}\|\|\Delta f + y\|/\sqrt{n_{\text{eff}}} > \lambda$, (see Appendix E.2)
560 are shaded gray. The nullclines $\Delta f_{t+1} = \Delta f_t$ and $\lambda_{t+1} = \lambda_t$ are shown as orange and white dashed
561 curves, respectively. Sharpness reduction, progressive sharpening, and divergent regions are colored
562 green, yellow, and blue. The gray arrows indicate the local vector field $\hat{G}(\Delta f, \lambda)$, which is the
563 direction of the updates. The training trajectories are depicted as black lines with arrows, with the
564 star marking the initialization. In all cases, $\eta_c = \sqrt{n_{\text{eff}}}/2$ (introduced in Section 5.2).

565 **Figure 3:** *UV model dynamics on the EoS manifold:*(left) Bifurcation diagram depicting late-time
566 limiting values of λ obtained by simulating Equation (3). (right) Bifurcation diagram of the UV
567 model. In both figures, $\|\mathbf{x}\| = 1, y = 2$ and $n_{\text{eff}} = 1$ and $\eta_c = 0.5$.

568 **Figure 4:** *Phase diagram of EoS:* (a) Heatmap of $\eta\bar{\lambda}^H/2$ of 3-layer ReLU FCNs with $s = 0$ trained
569 on a subset of 5,000 CIFAR-10 examples for 10k steps, with the weight variance σ_w^2 and learning
570 rate multiplier $c = \eta\lambda_0^H$ as axes. $\bar{\lambda}^H$ is obtained by averaging λ_t^H over last 200 steps. As the color
571 varies from blue to white, $\eta\bar{\lambda}^H/2$ increases, where the brightest white region indicates the EoS regime
572 with $\eta\bar{\lambda}^H/2 \geq 1$. (b) Same heatmap with fixed $\sigma_w^2 = 2.0$, but varying s continuously.

573 **Figure 5:** *EoS in synthetic vs real-datasets:* 2-layer linear FCN trained on (first row) 5,000 iid
574 random examples with unit output dimension and (second row) 5,000 CIFAR-10 examples. Different
575 columns correspond to the bifurcation diagram, late-time sharpness trajectories, and the power
576 spectrum of sharpness trajectories. Both models are trained for 10k steps using GD.

577 **Figure 10:** *Two-step phase portrait of UV model in $(\Delta f, \beta)$ phase plane:* These plots are equivalent
578 to Figure 2(d-f), but with training trajectory and local are plotted for every other step.

579 **Figure 14:** Sharpness and Weight Norm of 3-layer ReLU FCNs in SP with $\sigma_w^2 = 1/3$ and width
580 200, trained on a subset of CIFAR-10 with 5,000 examples using GD.

581 C.4 Sharpness measurement

582 We measure sharpness using the power iteration method with m iterations. Typically, $m = 20$
583 iterations ensure convergence. Exceptions requiring more iterates are discussed separately.

584 C.5 Power spectrum analysis

585 For a given signal $x'(t)$, we standardize the signal

$$x(t) = \frac{x'(t) - \mu}{\sigma}, \quad (10)$$

586 where μ is the mean and σ^2 is the variance of the signal. Subtracting the mean removes the zero
587 frequency component in the power spectrum. Next, consider the discrete Fourier transform $\mathcal{F}(\omega)$ of
588 $x(t)$:

$$\mathcal{F}(\omega) = \frac{1}{T} \sum_{t=0}^{T-1} e^{-i2\pi\omega t/T} x(t), \quad (11)$$

589 Then, the power spectrum is $P(\omega) = |\mathcal{F}(\omega)|^2$. The normalization by T in the Fourier transform
590 ensures that the sum of the power spectrum is equal to the variance of the signal, i.e., $\sum_{\omega} P(\omega) = \sigma^2$.

591 C.6 Estimation of Computational Resources Used

592 Most of our experiments, aside from the phase diagrams, required minimal computational resources,
593 estimated to be less than 50 A100 hours. In contrast, each phase diagram required 50 A100 hours,

594 totaling 500 A100 hours for all phase diagrams. Including initial experiments, we expect our total
 595 usage to be under 600 A100 hours.

596 D Sharpness dynamics of NTP networks

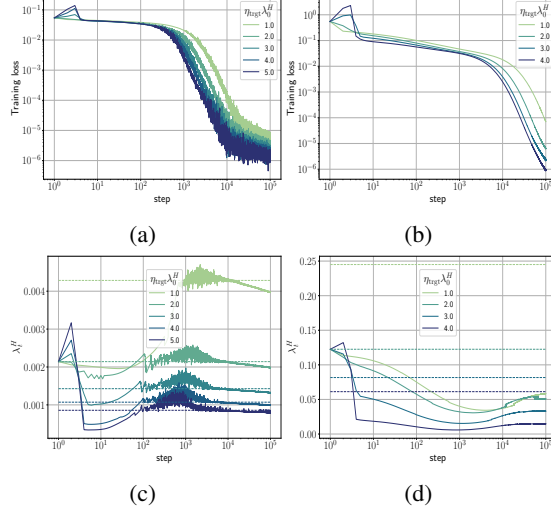


Figure 6: Training loss and sharpness trajectories of ReLU FCNs in NTP trained on a 5k subset of CIFAR-10 examples using MSE loss and GD: (a, c) $\sigma_w^2 = 0.5$, (b, d) $\sigma_w^2 = 2.0$.

597 Figure 6 shows that the sharpness dynamics of FCNs in NTP aligns with the behavior of FCNs in SP
 598 demonstrated in Figure 1.

599 E Properties of the UV model

600 E.1 Derivation of the Function Space Dynamics

601 Equations (1) and (2) can be derived using the gradient descent update equations:

$$U_{t+1} = U_t - \eta \frac{\Delta f_t \mathbf{v}_t \mathbf{x}^T}{\sqrt{n_{\text{eff}}}}, \quad (12)$$

$$\mathbf{v}_{t+1} = \mathbf{v}_t - \eta \frac{\Delta f_t U_t \mathbf{x}_t}{\sqrt{n_{\text{eff}}}}. \quad (13)$$

602 At step $t + 1$, the residual Δf_{t+1} can be written in terms of the gradient updates of U and \mathbf{v} :

$$\Delta f_{t+1} = f_{t+1} - y \quad (14)$$

$$= \frac{1}{\sqrt{n_{\text{eff}}}} \mathbf{v}_{t+1}^T U_{t+1} \mathbf{x} - y \quad (15)$$

$$= \frac{1}{\sqrt{n_{\text{eff}}}} \left(\mathbf{v}_t - \eta \frac{\Delta f_t U_t \mathbf{x}}{\sqrt{n_{\text{eff}}}} \right)^T \left(U_t - \eta \frac{\Delta f_t \mathbf{v}_t \mathbf{x}^T}{\sqrt{n_{\text{eff}}}} \right) \mathbf{x} - y \quad (16)$$

$$= \Delta f_t - \frac{\eta \Delta f_t}{n_{\text{eff}}} (\mathbf{x}^T U_t^T U_t \mathbf{x} + \mathbf{v}_t^T \mathbf{v}_t \mathbf{x}^T \mathbf{x}) + \frac{\eta^2 \|\mathbf{x}\|^2 \Delta f_t^2}{n_{\text{eff}}} \left(\frac{1}{\sqrt{n_{\text{eff}}}} \mathbf{x}^T U_t^T \mathbf{v}_t \right) \quad (17)$$

$$= \Delta f_t \left(1 - \eta \lambda_t + \frac{\eta^2 \|\mathbf{x}\|^2}{n_{\text{eff}}} \Delta f_t (\Delta f_t + y) \right). \quad (18)$$

603 Here, Δf_{t+1} only depends on Δf_t and λ_t . Similarly, we write down the λ_{t+1} using the gradient
 604 update equations:

$$\lambda_{t+1} = \frac{1}{n_{\text{eff}}} (\mathbf{x}^T U_{t+1}^T U_{t+1} \mathbf{x} + \mathbf{v}_{t+1}^T \mathbf{v}_{t+1} \mathbf{x}^T \mathbf{x}) \quad (19)$$

$$= \lambda_t - 4 \frac{\eta \|\mathbf{x}\|^2}{n_{\text{eff}}} \Delta f_t (\Delta f_t + y) + \frac{\eta^2 \|\mathbf{x}\|^2 \Delta f_t^2}{n_{\text{eff}}} \lambda_t \quad (20)$$

$$= \lambda_t + \frac{\eta \|\mathbf{x}\|^2}{n_{\text{eff}}} \Delta f_t^2 \left(\eta \lambda_t - 4 \frac{\Delta f_t + y}{\Delta f_t} \right). \quad (21)$$

605 Equations (18) and (21) form a closed system. This means that Δf_{t+1} and λ_{t+1} are completely
 606 described using Δf_t and λ_t . As a result, the complete dynamics of the UV model can be fully
 607 described using only these two variables with three parameters effective parameters η , $\frac{\|\mathbf{x}\|^2}{n_{\text{eff}}}$ and y .

608 E.2 Forbidden regions of the UV model

609 In this section, we utilize the non-negativity of λ to derive the condition for allowed regions within
 610 the phase plane for the UV model. Consider the function space equations written in terms of the
 611 pre-activation $\mathbf{h}(\mathbf{x}) = U\mathbf{x}$:

$$f(\mathbf{x}; \theta) = \frac{1}{\sqrt{n_{\text{eff}}}} \mathbf{v}^T \mathbf{h}(\mathbf{x}) \quad (22)$$

$$\lambda = \frac{1}{n_{\text{eff}}} (\|\mathbf{v}\|^2 \|\mathbf{x}\|^2 + \|\mathbf{h}\|^2). \quad (23)$$

612 Let $\cos(\mathbf{h}, \mathbf{v})$ denote the cosine similarity between \mathbf{v} and \mathbf{h} . Then, the network output is bounded as

$$|\cos(\mathbf{h}, \mathbf{v})| = \frac{\sqrt{n_{\text{eff}}} |\Delta f + y|}{\|\mathbf{v}\| \|\mathbf{h}\|} \leq 1 \quad (24)$$

613 Next, using $(\|\mathbf{v}\| \|\mathbf{x}\| - \|\mathbf{h}\|)^2 \geq 0$, we can bound the product $\|\mathbf{v}\| \|\mathbf{h}\|$ using λ

$$\frac{2\|\mathbf{x}\|}{\sqrt{n_{\text{eff}}}} |\Delta f + y| \leq \lambda. \quad (25)$$

614 The derived inequality describes the allowed phase plane regions for the UV model.

615 E.3 Fixed Points and Line

616 To identify the fixed points, we set $\Delta f_{t+1} = \Delta f_t$ and $\lambda_{t+1} = \lambda_t$ in Equations (1) and (2). This
 617 yields the dynamical fixed points of the UV model. Table 1 lists these fixed points along with their
 618 stability. Additionally, it provides the eigenvalues and eigenvectors of the Jacobian for the update
 619 maps described by Equations (1) and (2), evaluated at the fixed points.

620 E.4 The maximum learning rate η_{upper}

621 In Section 4, we stated that for $\eta > \eta_{\text{upper}} = 2\eta_c$, training diverges for all initializations except for
 622 those at the fixed points. Here, we justify this claim.

623 First, Figure 7 shows that as η approaches η_{upper} , fixed point III merges with fixed point II, reducing
 624 the convergence region to the EoS manifold. At this learning rate, the stability of fixed point II changes
 625 from saddle to unstable as the corresponding eigenvalue $(1 - \frac{\eta}{\eta_c})^2 \Big|_{\eta=2\eta_c}$ surpasses 1. Consequently,
 626 any initialization outside the EoS manifold results in divergence. Next, Figure 3(left) shows that on
 627 the EoS manifold training diverges for $\eta > \eta_{\text{upper}}$. This corroborates our initial claim.

	$(\Delta f^*, \lambda^*)$	eigenvalues	eigenvectors	Linear stability
I	$(0, \lambda)$ for $\lambda \geq \frac{2\ \mathbf{x}\ y}{\sqrt{n_{\text{eff}}}}$	$1, 1 - \eta\lambda$	$\begin{bmatrix} 0 \\ 1 \end{bmatrix}, \begin{bmatrix} \frac{n_{\text{eff}}\lambda}{4\ \mathbf{x}\ ^2 y} \\ 1 \end{bmatrix}$	$\begin{cases} \text{stable} & \eta\lambda < 2 \\ \text{unstable} & \eta\lambda > 2 \end{cases}$
II	$(-y, 0)$	$(1 - \frac{\eta\ \mathbf{x}\ y}{\sqrt{n_{\text{eff}}}})^2, (1 + \frac{\eta\ \mathbf{x}\ y}{\sqrt{n_{\text{eff}}}})^2$	$\begin{bmatrix} -\frac{\sqrt{n_{\text{eff}}}}{2\ \mathbf{x}\ } \\ 1 \end{bmatrix}, \begin{bmatrix} \frac{\sqrt{n_{\text{eff}}}}{2\ \mathbf{x}\ } \\ 1 \end{bmatrix}$	saddle
III	$(\frac{-2\sqrt{n_{\text{eff}}}}{\ \mathbf{x}\ \eta}, \frac{4}{\eta} - \frac{2\ \mathbf{x}\ y}{\sqrt{n_{\text{eff}}}})$	$9, 5 - \frac{2\eta\ \mathbf{x}\ y}{\sqrt{n_{\text{eff}}}}$	$\begin{bmatrix} \frac{n_{\text{eff}}}{\eta\ \mathbf{x}\ ^2 y} \\ 1 \end{bmatrix}, \begin{bmatrix} -\frac{\sqrt{n_{\text{eff}}}}{2\ \mathbf{x}\ } \\ 1 \end{bmatrix}$	unstable
IV	$(\frac{2\sqrt{n_{\text{eff}}}}{\ \mathbf{x}\ \eta}, \frac{4}{\eta} + \frac{2\ \mathbf{x}\ y}{\sqrt{n_{\text{eff}}}})$	$9, 5 + \frac{2\eta\ \mathbf{x}\ y}{\sqrt{n_{\text{eff}}}}$	$\begin{bmatrix} \frac{n_{\text{eff}}}{\eta\ \mathbf{x}\ ^2 y} \\ 1 \end{bmatrix}, \begin{bmatrix} \frac{\sqrt{n_{\text{eff}}}}{2\ \mathbf{x}\ } \\ 1 \end{bmatrix}$	unstable

Table 1: Fixed line (I) and points (II-IV) and corresponding eigenvalues and eigenvectors of the Jacobian of the update map in Equations (1) and (2). The stability is determined for $\eta < \eta_{\text{upper}} = 2\eta_c = 2\sqrt{n_{\text{eff}}}/\|\mathbf{x}\|y$. Above this threshold, training diverges for all initializations except for those at the fixed points, as demonstrated in Appendix E.4.

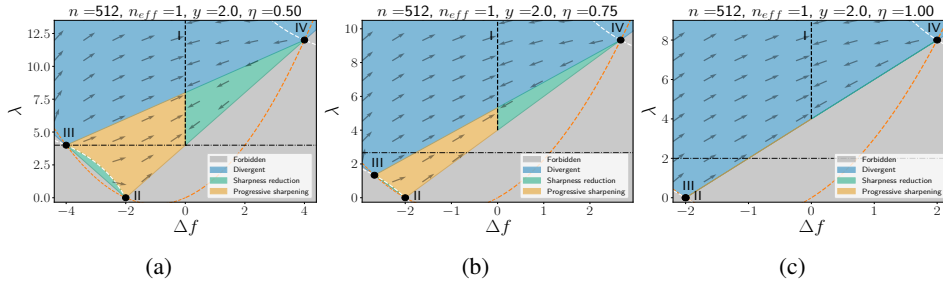


Figure 7: Phase portrait of the UV model for different learning rates η . The critical learning rate is $\eta_c = 0.5$ and the maximum learning rate is $\eta_{\text{upper}} = 1.0$.

628 E.5 Sharpness versus the trace of Hessian

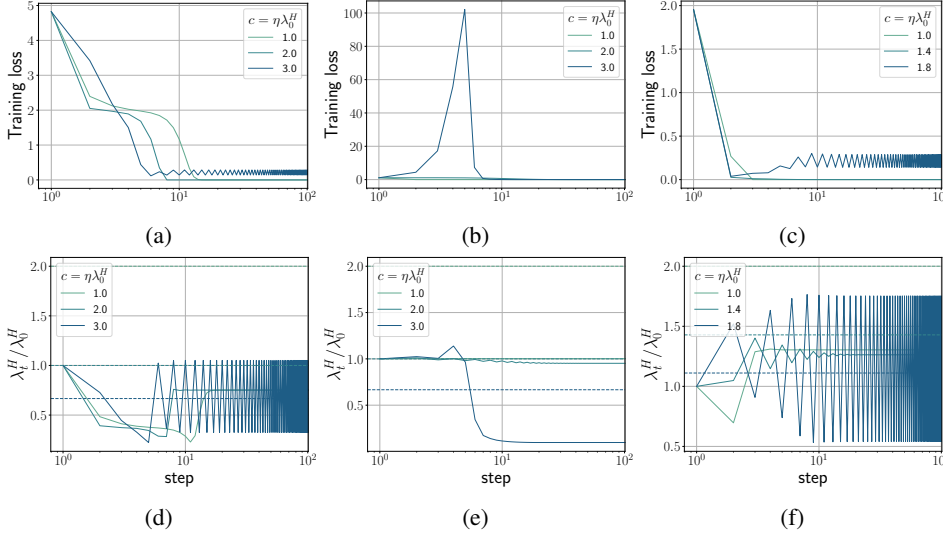


Figure 8: Training trajectories of the UV model trained on a single example with $\|\mathbf{x}\| = 1$ and $y = 2$ using MSE loss and GD: (a, d) NTP with $n = 1, \sigma_w^2 = 0.5$, (b, e) NTP with $n = 512, \sigma_w^2 = 1.0$, and (c, f) μP with $n = 512, \sigma_w^2 = 1.0$.

629 In this section, we show that the trace of the Hessian λ (which is also the scalar NTK in this case), is
630 an adequate proxy for sharpness. Figure 9 shows training trajectories of the UV model, with λ as a
631 proxy for sharpness and learning rate scaled as $\eta = k/\lambda_0$. These λ trajectories show similar trends to
632 those of λ^H observed in Figure 8, with one key difference: during early training, λ does not catapult

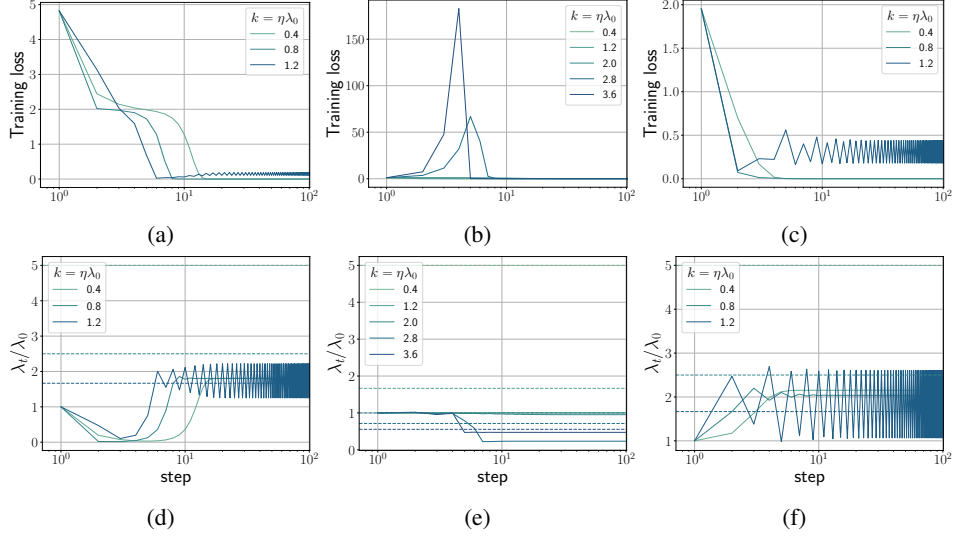


Figure 9: *UV model shows all four training regimes*: Training trajectories of the UV model trained on a single example (\mathbf{x}, y) with $\|\mathbf{x}\| = 1$ and $y = 2$ using MSE loss and gradient: (a, d) NTP with $n = 1$ and $\sigma_w^2 = 0.5$ (b, e) NTP with $n = 512$ and $\sigma_w^2 = 1.0$, and (c, f) μP with $n = 1$ and $\sigma_w^2 = 1.0$.

633 during early training at large widths (compare Figure 8(e) and Figure 9(e)). Otherwise, λ effectively
 634 captures other qualitative behavior of λ^H .

635 E.6 The distribution of residual and NTK at initialization

636 In this section, we compute the distribution of Δf and λ for the UV model at initialization. Consider
 637 the UV model written in terms of the pre-activation $h(x) = Ux$,

$$f(\mathbf{x}; \theta) = \frac{1}{\sqrt{n^{1-p}}} \mathbf{v}^T \mathbf{h}(\mathbf{x}) \quad (26)$$

$$\lambda = \frac{1}{n^{1-p}} (\|\mathbf{v}\|^2 \|\mathbf{x}\|^2 + \|\mathbf{h}\|^2), \quad (27)$$

638 with $v_i, U_{ij} \sim \mathcal{N}(0, \sigma_w^2/n^p)$. Then, each pre-activation h_i is normally distributed at initialization with
 639 zero mean and variance

$$\mathbb{E}_\theta[h_i^2] = \sum_{j,k=1}^{d_{\text{in}}} \langle U_{ij} U_{ik} \rangle x_j x_k = \sum_{j,k=1}^{d_{\text{in}}} \frac{\sigma_w^2}{n^p} \delta_{jk} x_j x_k = \frac{\sigma_w^2 \|\mathbf{x}\|^2}{n^p}. \quad (28)$$

640 Hence, each pre-activation is distributed as $h_i \sim \mathcal{N}(0, \sigma_w^2 \|\mathbf{x}\|^2/n^p)$. It follows that the network output
 641 is also normally distributed at initialization with zero mean and variance

$$\mathbb{E}_\theta[f_0^2] = \frac{1}{n^{1-p}} \sum_{i,j=1}^n \langle v_i v_j \rangle \langle h_i h_j \rangle = \frac{1}{n^{1-p}} \sum_{i=1}^n \frac{\sigma_w^2}{n^p} \frac{\sigma_w^2 \|\mathbf{x}\|^2}{n^p} = \frac{\sigma_w^4 \|\mathbf{x}\|^2}{n^p}. \quad (29)$$

642 Hence, the residual at initialization is distributed as $\Delta f_0 \sim \mathcal{N}(-y, \sigma_w^4 \|\mathbf{x}\|^2/n^p)$. Similarly, we can
 643 also compute the distribution of λ at initialization. The mean value of λ is given by

$$\mathbb{E}_\theta[\lambda_0] = \frac{1}{n^{1-p}} (\|\mathbf{x}\|^2 \langle \|\mathbf{v}\|^2 \rangle + \langle \|\mathbf{h}\|^2 \rangle) = 2\sigma_w^2 \|\mathbf{x}\|^2, \quad (30)$$

644 where we have used $\langle \|\mathbf{v}\|^2 \rangle = \sigma_w^2 n^{1-p}$ and $\langle \|\mathbf{h}\|^2 \rangle = \sigma_w^2 \|\mathbf{x}\|^2 n^{1-p}$. Using similar computations,
 645 the second moment of λ is given by:

$$\mathbb{E}_\theta[\lambda_0^2] = \frac{1}{n^{2-2p}} (\|\mathbf{x}\|^4 \langle \|\mathbf{v}\|^4 \rangle + \langle \|\mathbf{h}\|^4 \rangle + 2\|\mathbf{x}\|^2 \langle \|\mathbf{v}\|^2 \|\mathbf{h}\|^2 \rangle) = \frac{4(n+1)}{n} \sigma_w^4 \|\mathbf{x}\|^4. \quad (31)$$

646 Hence, the λ at initialization is distributed as $\lambda_0 \sim \mathcal{N}(2\sigma_w^2 \|\mathbf{x}\|^2, 4\sigma_w^4 \|\mathbf{x}\|^4/n)$.

647 E.7 EoS manifold is a dynamical attractor

648 To demonstrate that late time trajectories for $\eta > \eta_c$ converge to the EoS manifold, we define
 649 $\beta := \frac{\sqrt{n_{\text{eff}}}}{2\|\mathbf{x}\|} \lambda - (\Delta f + y)$. β lies on the direction orthogonal to the EoS manifold, such that $\beta = 0$
 650 corresponds to the manifold itself, while $\beta < 0$ is forbidden. Under this transformation, β updates as
 651 $\beta_{t+1} = \beta_t (1 + \frac{\eta \|\mathbf{x}\| \Delta f_t}{\sqrt{n_{\text{eff}}}})^2$. It follows that $\beta^* = 0$ stays invariant under the dynamics and defines a
 652 nullcline.

653 Due to oscillations in Δf near convergence, it is instructive to examine the two-step dynamics [1, 5],
 654 compactly denoted as $(\Delta f_{t+2}, \lambda_{t+2}) := M^{(2)}(\Delta f_t, \lambda_t)$. Figure 10 shows the two-step trajectories
 655 and the corresponding vector field $\hat{G}^{(2)}(\Delta f, \beta)$ in the $(\Delta f, \beta)$ plane.

656 We observe that there exists a critical η_c such that for $\eta < \eta_c$, $\hat{G}^{(2)}(\Delta f, \beta)$ points towards the
 657 stable zero-loss line (see Figure 10(a)). By comparison, for $\eta > \eta_c$, all points along the zero-loss
 658 line become unstable and the vector field directs towards points on the $\beta = 0$ line, as shown in
 659 Figure 10(b). The critical η_c for which all points on the zero-loss line become unstable thus gives a
 660 necessary condition for EoS:

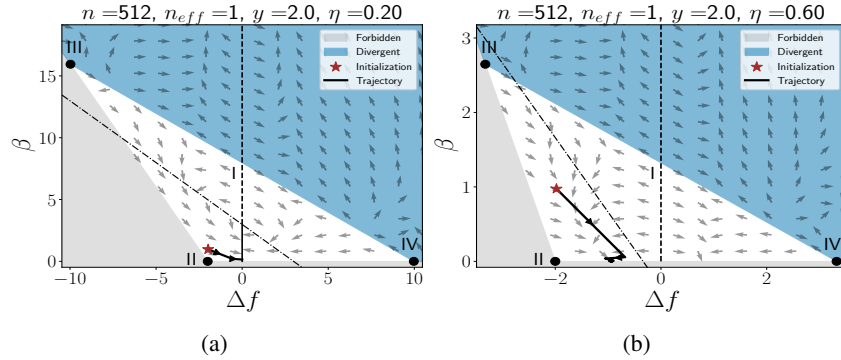


Figure 10: These plots are equivalent to Figure 2(a, d) ($\eta_c = 0.5$), but with training trajectory and local vector field plotted for every other step in $(\Delta f, \beta)$ plane. The tilted dash-dotted line indicates the $\eta\lambda = 2$ line.

661 E.8 Critical learning rate for edge of stability

662 In this section, we estimate the required condition on the learning rate for the UV model to exhibit
 663 EoS. We specifically focus on the case with $y > 0$ as for $y = 0$, λ can only decrease. As a result,
 664 the model does not exhibit progressive sharpening and EoS. In Section 5, we observed that the EoS
 665 occurs as the zero-loss minima with the smallest λ becomes unstable. From Equation (25) it follows
 666 that the smallest λ with zero loss is

$$\lambda_{\min} = \frac{2\|\mathbf{x}\|y}{\sqrt{n_{\text{eff}}}}. \quad (32)$$

667 This minimum becomes unstable if the learning rate η exceeds a critical value η_c , given by

$$\eta_c = \frac{\sqrt{n_{\text{eff}}}}{\|\mathbf{x}\|y} \quad (33)$$

668 It is worth noting that this is a necessary condition for λ to oscillate around $2/\eta$. Otherwise, training
669 converges to the zero-loss minimum with $\lambda = \lambda_{\min}$ for $\eta < \eta_{\max}$.

670 We can also derive the exact same result by analyzing the dynamics on the EoS manifold. As
671 discussed in Section 5.2, the dynamics on the EoS manifold is given by the map $\Delta f_{t+1} = M(\Delta f_t)$,
672 where

$$M(\Delta f) = \Delta f \left(1 - \frac{2\eta\|\mathbf{x}\|}{\sqrt{n_{\text{eff}}}}(\Delta f + y) + \left(\frac{\eta\|\mathbf{x}\|}{\sqrt{n_{\text{eff}}}} \right)^2 \Delta f(\Delta f + y) \right). \quad (34)$$

673 As demonstrated in Section 5.2, EoS in the UV model follows the period doubling route to chaos,
674 with the period two cycle marking the onset. Hence, the conditions required for emergence of the
675 period two cycle are also the necessary conditions for EoS. Consider the two-step dynamics on the
676 EoS manifold given by the map $M^2(\Delta f) := M(M(\Delta f))$. This map has six fixed points (excluding
677 three fixed points of the map M) summarized below

$$\Delta f^* = \frac{\eta\tilde{x}(1 - \eta\tilde{x}y) \pm \sqrt{\eta^2\tilde{x}^2(\eta\tilde{x}y - 1)(3 + \eta\tilde{x}y)}}{2\eta^2\tilde{x}^2} \quad (35)$$

$$\Delta f^* = \frac{3 + h(\eta, \tilde{x}, y) \pm \eta\tilde{x}(\pm y + \sqrt{2}\sqrt{-\frac{-5 + h(\eta, \tilde{x}, y) + \eta\tilde{x}y(-2 - \eta\tilde{x}y + h(\eta, \tilde{x}, y))}{\eta^2\tilde{x}^2}})}{4\eta\tilde{x}}. \quad (36)$$

678 Here $\tilde{x} = \frac{\|\mathbf{x}\|}{\sqrt{n_{\text{eff}}}}$ and $h(\eta, \tilde{x}, y) = \sqrt{-7 + \tilde{x}y\eta(2 + \tilde{x}y\eta)}$. For the fixed points to exist, we require the
679 expressions inside the square root to be non-negative, i.e.,

$$\left(\frac{\eta\|\mathbf{x}\|y}{\sqrt{n_{\text{eff}}}} - 1 \right) \left(\frac{\eta\|\mathbf{x}\|y}{\sqrt{n_{\text{eff}}}} + 3 \right) \geq 0 \implies \eta \geq \eta_1 = \frac{\sqrt{n_{\text{eff}}}}{\|\mathbf{x}\|y} \quad (37)$$

$$\frac{\eta\|\mathbf{x}\|y}{\sqrt{n_{\text{eff}}}} \left(\frac{\eta\|\mathbf{x}\|y}{\sqrt{n_{\text{eff}}}} + 2 \right) - 7 \geq 0 \implies \eta \geq \eta_2 = \frac{(\sqrt{32} - 2)}{2} \frac{\sqrt{n_{\text{eff}}}}{\|\mathbf{x}\|y}. \quad (38)$$

680 As $\eta_1 < \eta_2$, the necessary condition for the period two cycle to emerge is $\eta > \eta_1 = \sqrt{n_{\text{eff}}}/\|\mathbf{x}\|y$, which
681 coincides with the condition obtained earlier in this section.

682 F The effect of batch size on the four training regimes

683 In this section, we examine the effect of batch size B on the results presented in the main text. We
684 find that our conclusions are robust for reasonable batch sizes around $B \approx 512$. For even smaller
685 batch sizes, the dynamics becomes noise-dominated, and separating the inherent dynamics from noise
686 becomes challenging. This observation further supports the use of SGD to reduce the computational
687 cost of experiments in the subsequent sections involving CNNs and ResNets.

688 Figure 11 shows that SGD trajectories of FCNs in SP begin to deviate from their GD counterpart
689 significantly for batch sizes around $B \approx 128$. In contrast, for μP networks this deviation begins at a
690 larger batch size of $B \approx 512$ as shown in Figure 12. Figure 13 show training trajectories of CNNs
691 and ResNets trained SGD with batch size $B = 512$. These results further exemplify that four regimes
692 of training are generically observed for reasonable batch sizes.

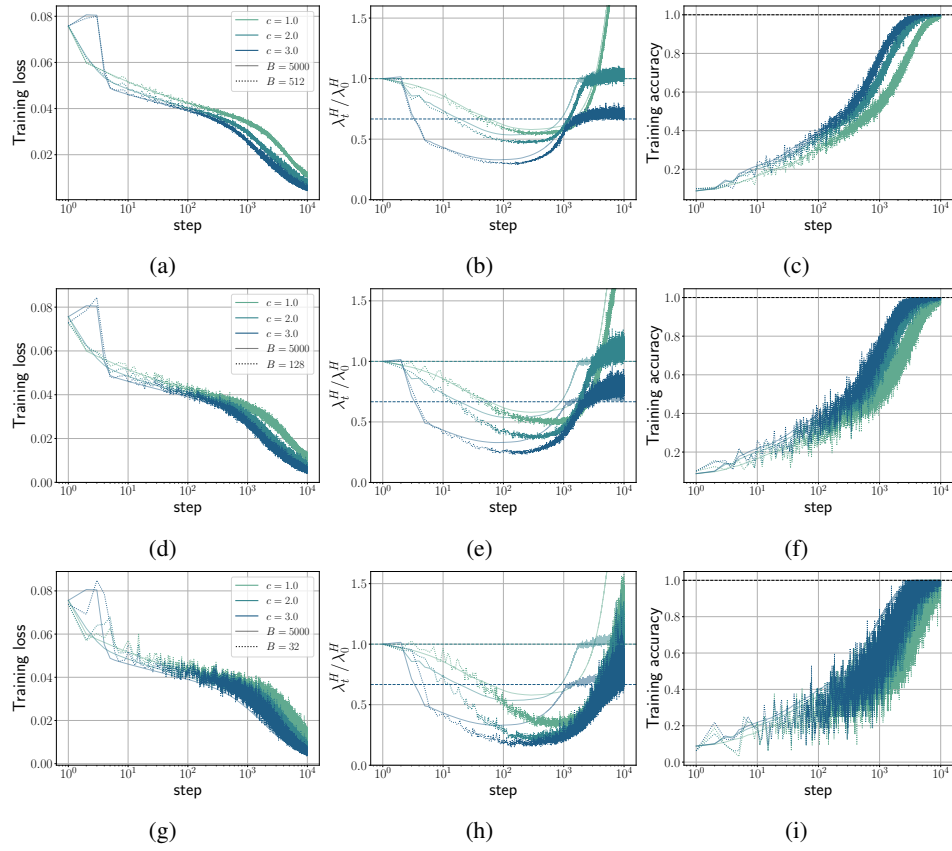


Figure 11: Comparison of SGD trajectories with their GD counterpart for a 3-layer FCNs in SP with $\sigma_w^2 = 0.5$ trained on a subset of CIFAR-10 consisting of 5,000 training examples with MSE loss. The learning rate is scaled as $\eta = c/\lambda_0^H$ and batch sizes (a-c) $B = 512$, (d-f) $B = 128$, (g-i) $B = 32$ are considered. GD trajectories are plotted using solid lines with transparency.

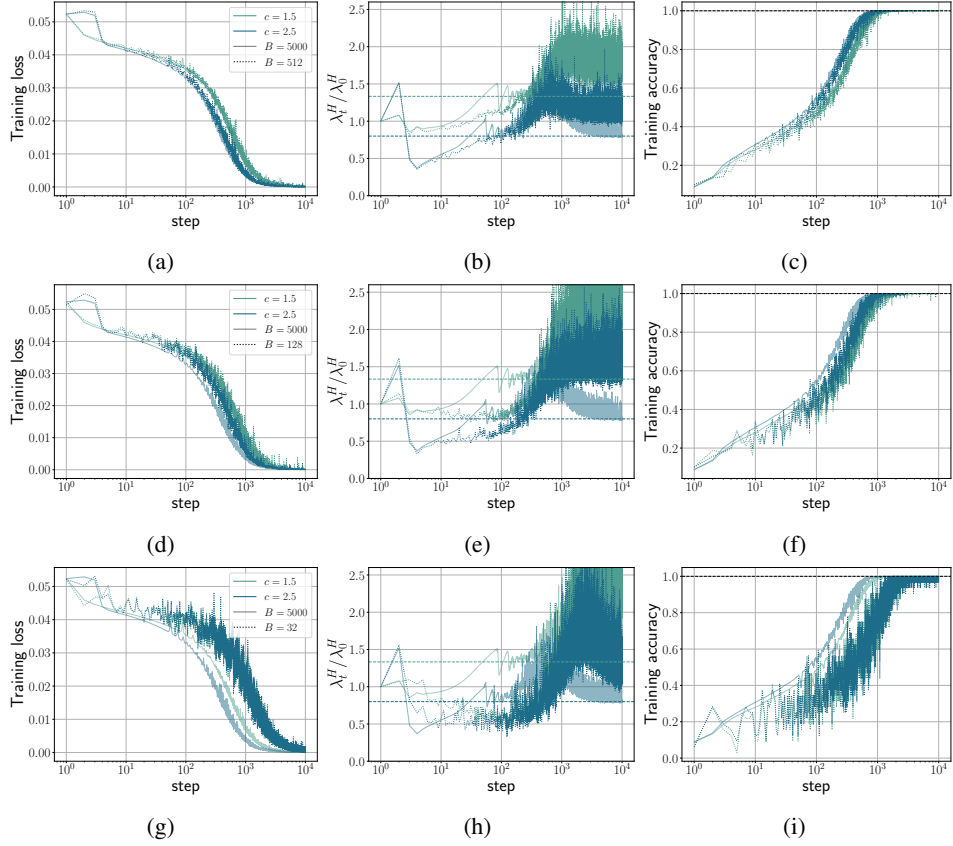


Figure 12: Same setting as Figure 11, except we used 3-layer FCNs in μP with $\sigma_w^2 = 2.0$.

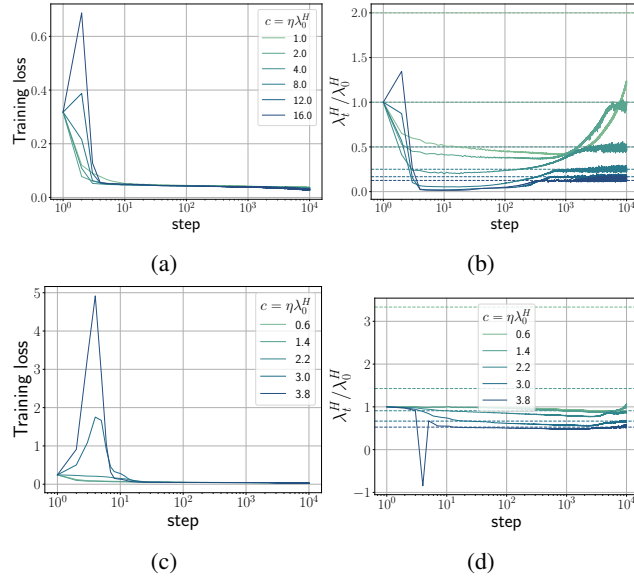


Figure 13: Training trajectories of (a, b) a 5-layer CNN in SP with $n = 64$, and (c, d) ResNet-18 with LayerNorm in SP, also with $n = 64$. Both models are trained on the CIFAR-10 dataset with MSE loss using SGD. The learning rate is scaled as $\eta = c/\lambda_0^H$ and batch size is $B = 512$. In panel (d), λ_t^H becomes negative during early training. This is due to the power iteration method returning the largest eigenvalue by magnitude.

693 **G Sharpness-weight norm correlation**

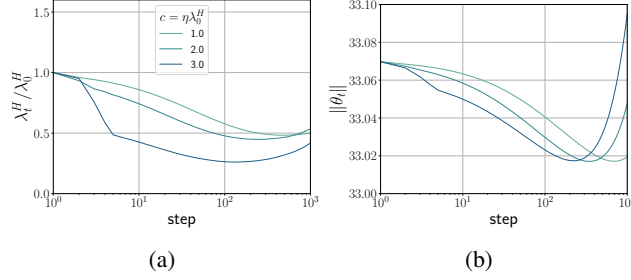


Figure 14: Sharpness and Weight Norm of 3-layer ReLU FCNs in SP with $\sigma_w^2 = 1/3$, trained on a subset of CIFAR-10 with 5,000 examples using GD.

694 Section 5.1 reveals that several aspects of the training dynamics are controlled by the fact that, for a
 695 wide variety of initializations, at early times trajectories move closer to the saddle point II, resulting
 696 in an interim decrease in λ (also proportional to the weight norm in this case), before eventually
 697 increasing. This critical point where all parameters are zero also exists in real-world models. We
 698 thus anticipate that in real-world models, the origin of the four training regimes may be related to a
 699 similar mechanism. This would predict a decrease in weight norm as training passes near the saddle
 700 point, followed by an eventual increase.

701 Figure 14 validates this hypothesis. During the sharpness reduction and intermediate saturation
 702 regimes, we see a decrease in the weight norm, followed by an increase in the weight norm as the
 703 network undergoes progressive sharpening, following the prediction from the UV model. Similar
 704 correlations between the last layer weight norm and sharpness are utilized by [32] to analyze the
 705 EoS phase. By comparison, we focus on the correlation between sharpness and weight norm during
 706 early training to attribute the emergence of four regimes to the critical point corresponding to all
 707 parameters being zero. In Appendix G, we provide further evidence for this correlation between
 708 sharpness and weight norm, extending this relationship to CNNs and ResNets.

709 This section presents additional results for Appendix A, further supporting the relationship between
 710 sharpness and weight norm during training.

711 Figure 15 is an extended version of Figure 14, where we plotted the whole training trajectories and
 712 measured Pearson correlation

$$\text{Cor}(\|\theta_t\|, \lambda_t^H / \lambda_0^H) := \frac{\sum_{t'=1}^t (\theta_{t'} - \bar{\theta}_t) \left(\lambda_{t'}^H / \lambda_0^H - \overline{(\lambda^H / \lambda_0^H)}_t \right)}{\sqrt{\sum_{t'=1}^t (\theta_{t'} - \bar{\theta}_t)^2 \sum_{t'=1}^t \left(\lambda_{t'}^H / \lambda_0^H - \overline{(\lambda^H / \lambda_0^H)}_t \right)^2}} \quad (39)$$

713 Here $t \geq 2$ and $\bar{\theta}_t = (\sum_{t'=1}^t \theta_{t'}) / t$.

714 Figure 16 shows the weight norm of each layer separately for the experiment in Figure 14. This result
 715 shows a high correlation between weight norm and sharpness through training.

716 We also confirm these correlations between weight norm and sharpness in CNNs for the experiment
 717 in Figure 13(a, b).

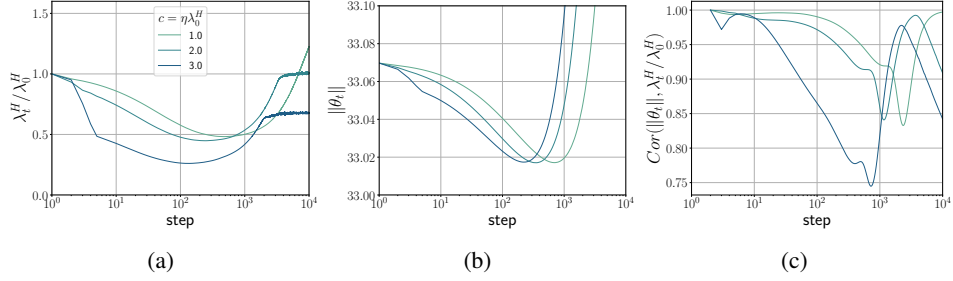


Figure 15: Sharpness and Weight Norm of 3-layer ReLU FCNs in SP with $\sigma_w^2 = 1/3$, trained on a subset of CIFAR-10 with 5,000 examples using GD.

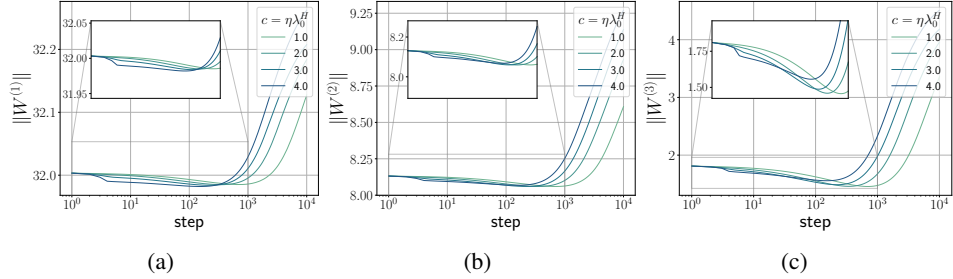


Figure 16: Weight Norm of each layer in 3-layer ReLU FCNs (same experiments as Figure 14): (a, b, c) SP with $\sigma_w^2 = 1/3$. All results are obtained by training on a subset of CIFAR-10 with 5,000 examples using GD.

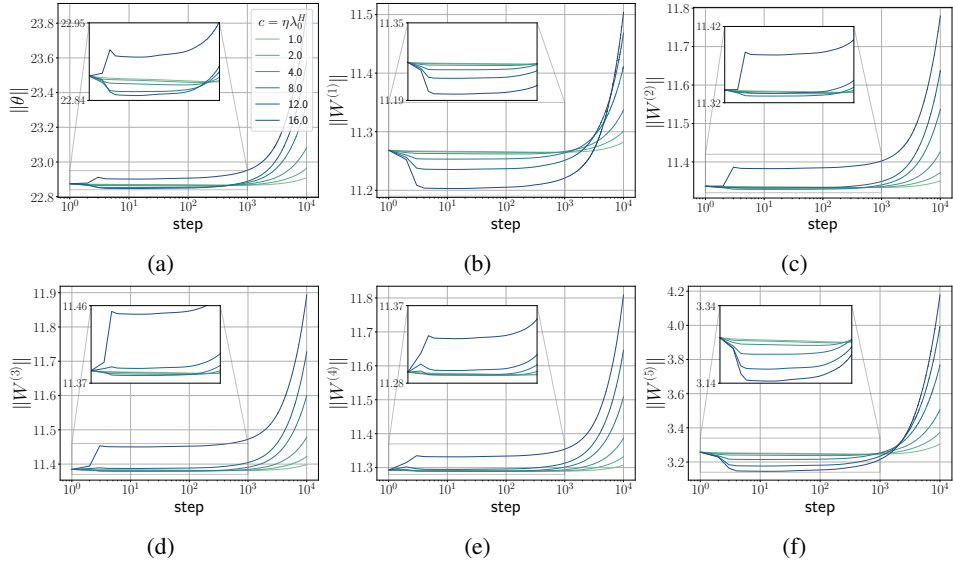


Figure 17: Weight Norm of each layer for 5-layer CNNs in SP (same experiments as Figure 13(a, b)): (a) Total weight norm; (b-f) Weight norm of each layer. We see that for $c = 16$, the initial catapult in sharpness λ^H (Figure 13(b)) is accompanied by a catapult in total weight norm. Notably, the total weight norm and per-layer weight norm, whether catapults (a, c-e) or not (b, f), show a decreasing trend during the early sharpness decreasing stage, followed by an eventual increase.

718 H Additional Phase diagrams of EoS

719 This section demonstrates additional phase diagrams of EoS and quantifies the effect of batch size in
 720 the EoS regime. Figure 18 shows phase diagrams of EoS for FCNs trained on CIFAR-10 with MSE

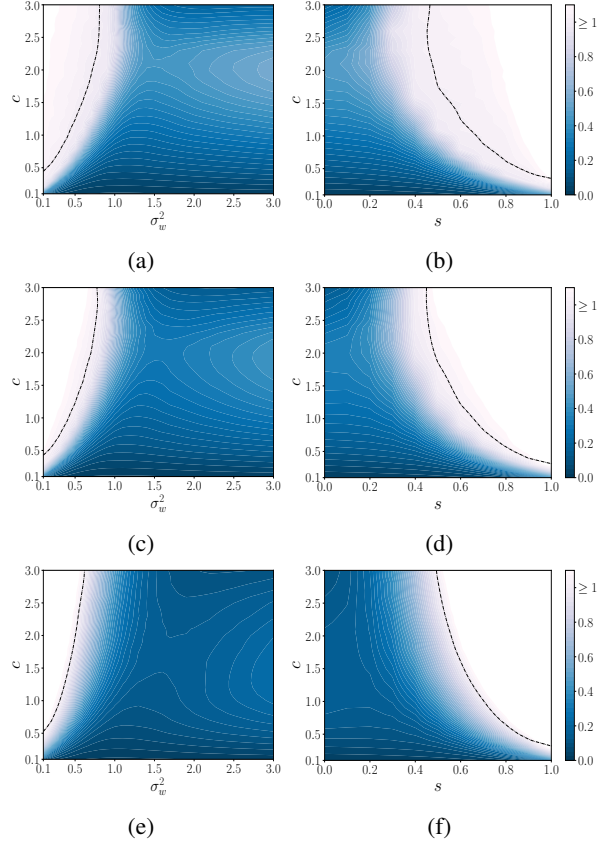


Figure 18: Phase diagram of EoS for 3-layer FCNs trained on CIFAR-10 with MSE loss using SGD with three different batch sizes: (a, b) $B = 512$, (c, d), $B = 128$, and (e, f) $B = 32$. The color indicates the value of $\eta\bar{\lambda}^H/2$, where $\bar{\lambda}^H$ is obtained by averaging λ_t^H over the last 200 steps. Except for the batch size, all settings are identical to Figure 4. Black dash-dotted lines indicate the phase boundary $\eta\bar{\lambda}^H/2 = 1$. For clarity, these lines are generated from data smoothed with a Gaussian kernel.

721 loss using SGD for 10,000 steps for three different batch sizes. We observe that as the batch size
 722 decreases, λ_t^H oscillates at a value different from $2/\eta$ depending on σ_w^2 and s . For large σ_w^2 and small
 723 s , λ^H favors a smaller value for smaller batch size, which is in agreement with the observation in [7].
 724 In contrast, λ^H can be larger than $2/\eta$ for small σ_w^2 and large s at late training times.

725 Figures 19 and 20 show the phase diagrams of EoS for CNNs and ResNets trained on the CIFAR-10
 726 dataset with MSE loss using SGD for 10,000 steps with learning rate $\eta = c/\lambda_0^H$ and batch size
 727 $B = 512$. In contrast to the FCN phase diagrams, these architectures exhibit EoS behavior at smaller
 728 values of s and larger values of σ_w^2 , indicating their implicit bias towards EoS. Moreover, we observe
 729 in ResNets, that EoS is less sensitive to change σ_w^2 , likely due to a combination of LayerNorm and
 730 residual connections [10].

731 It is worth noting that EoS boundaries in these phase diagrams are time-dependent. For instance,
 732 models close to the EoS boundary may eventually reach EoS on training longer (see Figure 13(b)
 733 $c = 1.0$ for example), causing a shift in the EoS boundary. Nevertheless, models with small learning
 734 rates, large σ_w^2 , and small s may never show EoS behavior, regardless of training duration, as predicted
 735 by the UV model and seen in Figure 1(e).

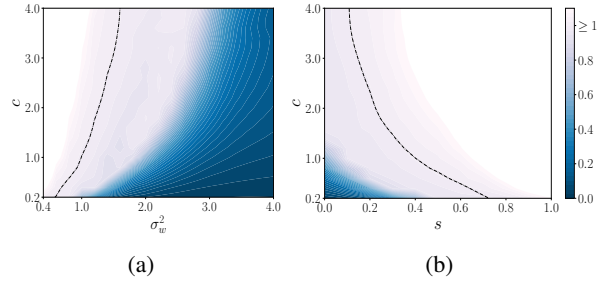


Figure 19: Phase diagram of EoS for 5-layer CNNs in SP with width $n = 64$ trained with MSE loss using SGD for 10,000 steps with learning rate $\eta = c/\lambda_0^H$ and batch size $B = 512$.

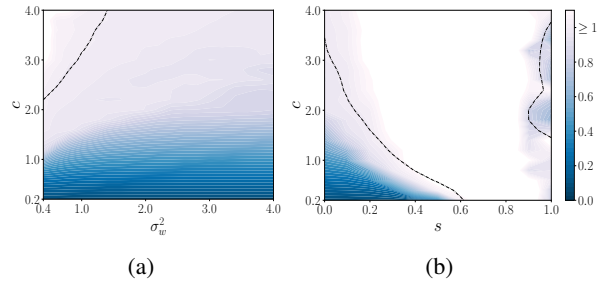


Figure 20: Phase diagram of EoS for ResNet-18 in SP with width $n = 64$ trained with MSE loss using SGD for 10,000 steps with learning rate $\eta = c/\lambda_0^H$ and batch size $B = 512$. For $s = 1$, the average eigenvalue $\bar{\lambda}^H$ is observed to be less than $2/\eta$. Upon detailed investigation of the trajectories, we found that λ_t^H oscillates around a value lower than $2/\eta$. We leave this anomalous behavior as an observation.

736 **I Route to chaos**

737 **I.1 Route to EoS in real datasets**

738 This section presents additional bifurcation diagrams for different architectures and datasets. Fig-
 739 ures 21 to 23 show the bifurcation diagrams, sharpness trajectories and the associated power spectrum
 740 of 4-layer ReLU FCNs in SP trained on MNIST, Fashion-MNIST and CIFAR-10 datasets with MSE
 741 loss using GD. Similarly, Figures 24 and 25 show these results for CNNs and ResNets trained on
 742 CIFAR-10 with MSE loss using SGD with batch size $B = 512$. These results show the reminiscent
 743 of the period-doubling route to chaos observed in different architectures and datasets. In all figures,
 744 we choose the smallest and largest learning rate exhibiting EoS for plotting the trajectories and power
 745 spectrum. The structured route to chaos in realistic experiments can be disrupted due to a variety of
 746 reasons. Below, we discuss a few of them.

747 **Measurement of only the top eigenvalue of Hessian:** In our experiments, we only measured the
 748 top eigenvalue of the Hessian. However, when multiple eigenvalues of Hessian enter EoS, plotting
 749 only the top eigenvalue of Hessian is a projection that could obscure all the structured routes to chaos
 750 that the system may exhibit.

751 **The effect of correlations in real-world datasets:** Real-world datasets inherently contain
 752 correlations between different samples (x, y) . These correlations can be quantified using the
 753 input-input covariance matrix $\Sigma_{XX} = XX^T \in \mathbb{R}^{d_{in} \times d_{in}}$ and output-input covariance matrices
 754 $\Sigma_{YX} = YX^T \in \mathbb{R}^{d_{out} \times d_{in}}$. In Appendix I.2, we find that a key determining factor in observing
 755 route-to-chaos is whether the power spectrum of Σ_{XX} is flat or exhibits power law decay. We show
 756 that power-law decay in the singular values of the Σ_{XX} results in long-range correlations in time and
 757 dense sharpness bands observed in real datasets.

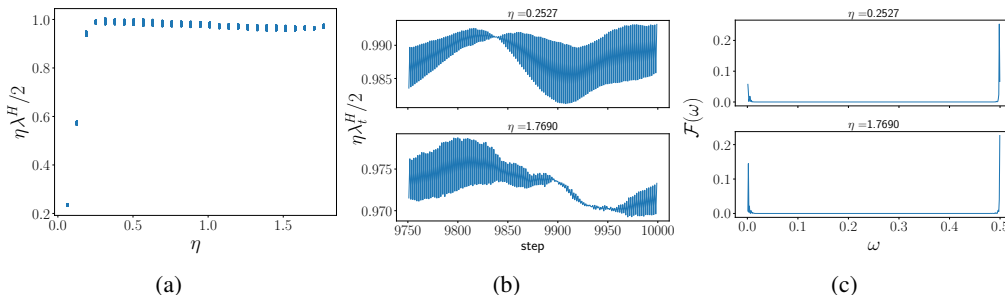


Figure 21: 4-layer FCN in SP with width $n = 512$ trained on a subset of 5000 examples of MNIST with MSE loss using GD. Both power spectrums are computed using the last 1000 steps of the corresponding trajectories.

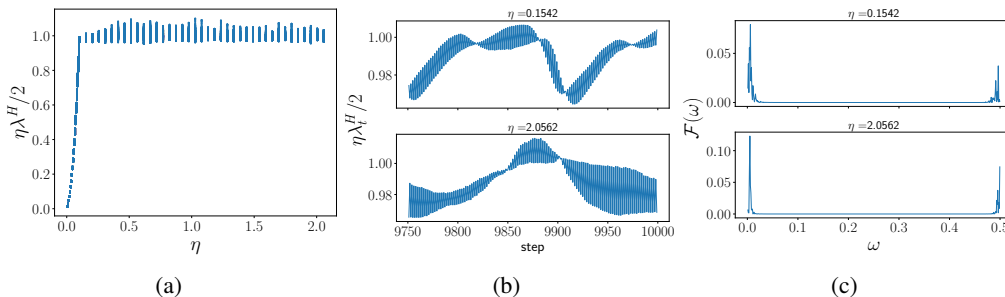


Figure 22: 4-layer FCN in SP with width $n = 512$ trained on a subset of 5000 examples of Fashion-MNIST with MSE loss using GD. Both power spectrums are computed using the last 1000 steps of the corresponding trajectories.

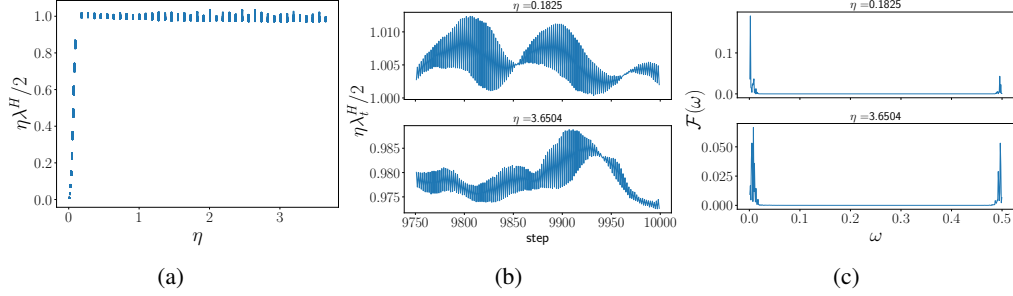


Figure 23: 4-layer FCN in SP with width $n = 512$ trained on a subset of 5000 examples of CIFAR-10 with MSE loss using GD. Both power spectrums are computed using the last 1000 steps of the corresponding trajectories.

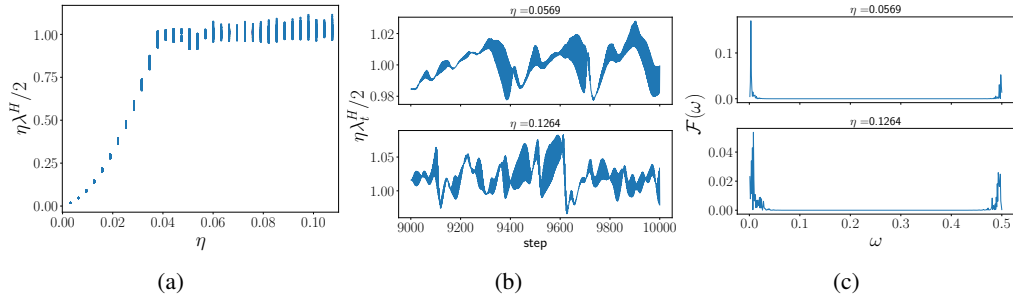


Figure 24: 5-layer CNN in SP with width $n = 32$ trained on a subset of 1000 examples of CIFAR-10 with MSE loss using GD.

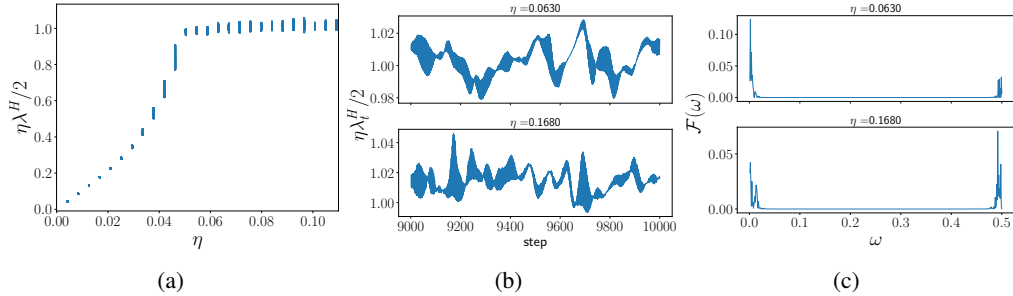


Figure 25: ResNet-18 in SP with width $n = 32$ trained on a subset of 1000 examples of CIFAR-10 with MSE loss using GD.

758 **I.2 The effect of power-law trends in data on sharpness trajectories**

759 In this section, we analyze a 2-layer linear FCN trained on the power law dataset described in
 760 Appendix C.1 to understand the origin of long-range correlations in sharpness trajectories and dense
 761 sharpness bands in realistic datasets.

762 Figure 26 shows the bifurcation diagram, late time trajectories, and the associated power spectrum
 763 of the network trained on the power-law dataset with the same $A_x = 1.0$ and $A_y = 1.0$, for four
 764 different combinations of power-law exponents: (i) $B_x = 0.0, B_y = 0.0$, (ii) $B_x = 1.0, B_y = 0.0$,
 765 (iii) $B_x = 0.0, B_y = 1.0$, and (iv) $B_x = 1.0, B_y = 1.0$. We observe that a power-law trend to the
 766 singular values of the input matrix results in dense sharpness bands observed in real datasets. It is
 767 worth noting that this is one way to obtain dense sharpness bands and in general, there can be many
 768 other methods.

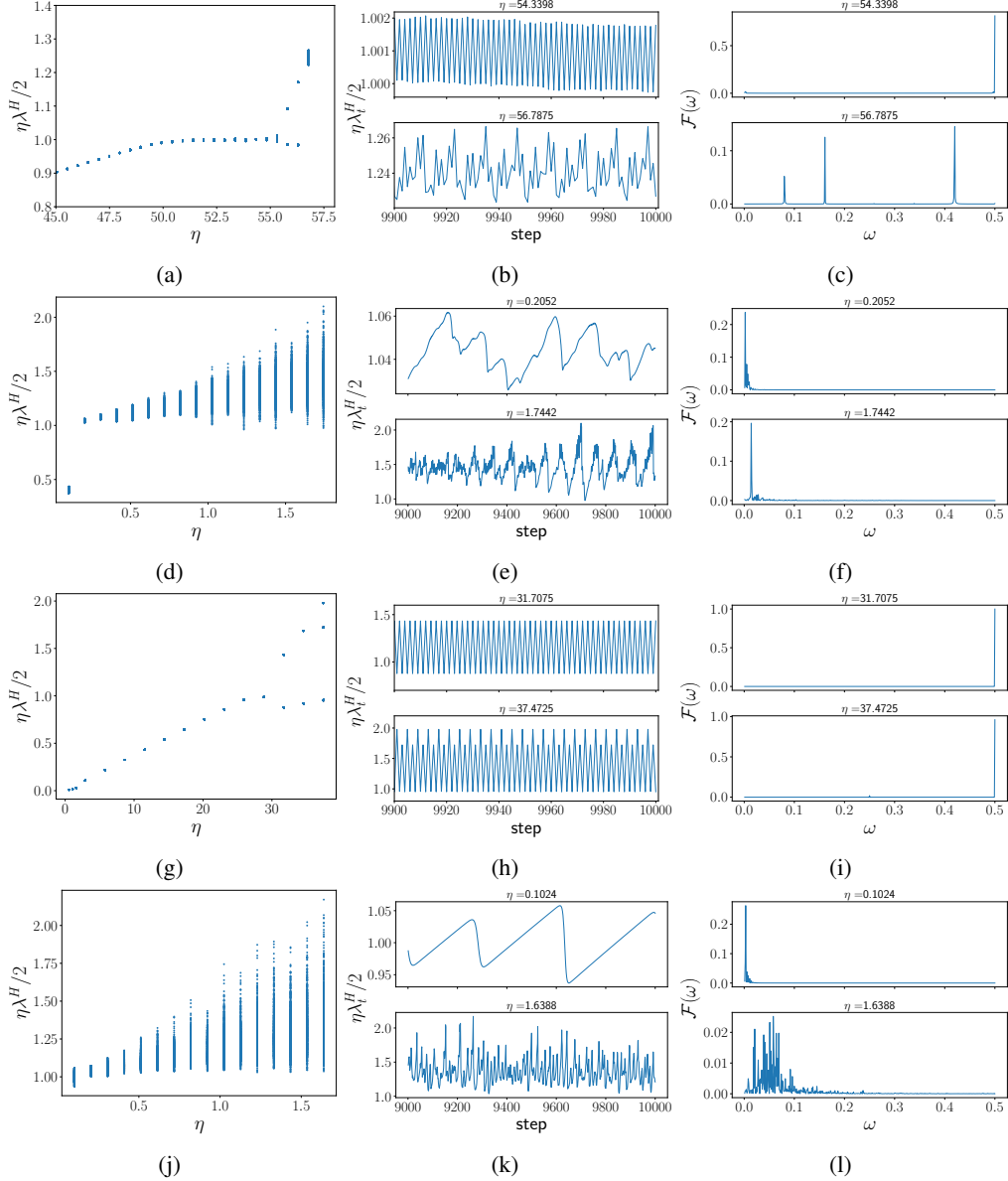


Figure 26: Bifurcation diagrams, late-time sharpness trajectories, and power spectrum of a 2-layer linear network trained on the power-law dataset for different parameter values: (a-c) $B_x = 0.0, B_y = 0.0$, (d-f) $B_x = 1.0, B_y = 0.0$, (g-i) $B_x = 0.0, B_y = 1.0$, and (j-l) $B_x = 1.0, B_y = 1.0$. All power spectrums are computed using the last 1000 steps of the corresponding trajectories.

769 **I.3 Route to chaos in synthetic datasets**

770 In this section, we analyze the route to chaos in synthetic datasets to gain insights into the dense
 771 sharpness bands in realistic datasets. We considered two datasets, defined as follows:

772 **Teacher-student dataset:** Consider a teacher FCN $f : \mathbb{R}^{d_{\text{in}}} \rightarrow \mathbb{R}^{d_{\text{out}}}$ with $d_{\text{in}} = 3072, d_{\text{out}} = 10$,
 773 depth d , and width $n = 512$ in Standard Parameterization. Then, we construct a teacher-student
 774 dataset (X, Y) consisting of $P = 5000$ examples with $\mathbf{x}^\mu \sim \mathcal{N}(0, I)$ and $\mathbf{y}^\mu = \mathbf{f}(\mathbf{x}^\mu; \theta_0)$. Next,
 775 we train a student FCN with the same depth d and depth n as the teacher FCN on this dataset.

776 Figures 27 and 28 show the bifurcation diagram, late time sharpness trajectories and the associated
 777 power spectrum of linear and ReLU FCNs trained on the teacher-student task. These figures show

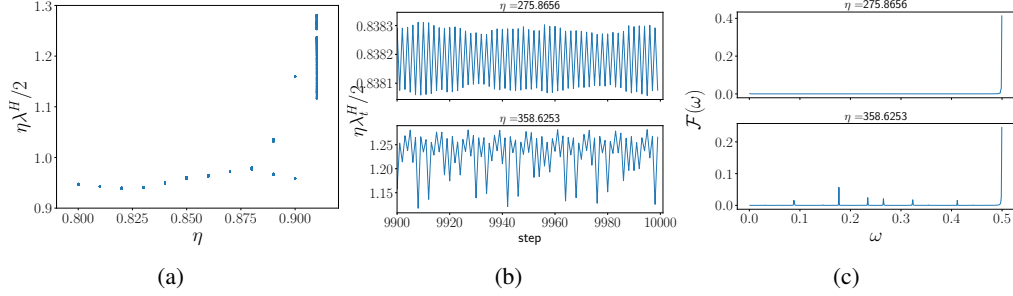


Figure 27: 2-layer linear FCN in μP trained on the teacher-student task. Both power spectrums are computed using the last 1000 steps of the corresponding trajectories.

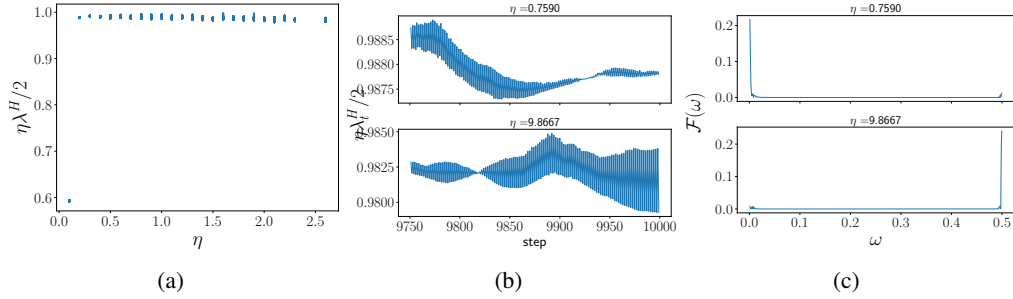


Figure 28: 4-layer ReLU FCNs in μP trained on the teacher-student task. Both power spectrums are computed using the last 1000 steps of the corresponding trajectories.

778 that while linear FCN shows the period-doubling route to chaos, ReLU FCN shows long-range
 779 correlations as observed in real datasets.

780 **Generative dataset:** Consider a 5-layer CNN $f(x, \theta)$ in SP with $n = 64$, trained on the CIFAR-10
 781 dataset with MSE loss using SGD with learning rate $\eta = 12/\lambda_0^H$ and momentum $m = 0.9$ for 100k
 782 steps. This model achieves a test accuracy of 76.9%. Then, we construct a generative image dataset
 783 (X, Y) consisting of $P = 5000$ examples with $x^\mu \sim \mathcal{N}(0, I)$ and $y^\mu = f(x^\mu; \theta)$. Next, we train an
 784 FCN in SP with depth d , width n , and weight variance $\sigma_w^2 = 0.5$ on the generated dataset.

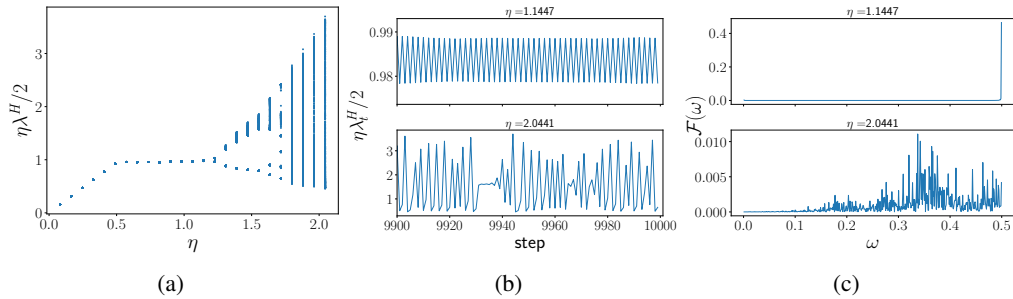


Figure 29: 4-layer linear FCNs in SP with $\sigma_w^2 = 0.5$ trained on the generative CIFAR-10 task with MSE loss using GD. Both power spectrums are computed using the last 1000 steps of the corresponding trajectories.

785 Figures 29 and 30 show the bifurcation diagram, late time trajectories and the associated power
 786 spectrum of a 4-layer ReLU FCN with linear and ReLU activations, trained on the generative CIFAR-
 787 10 dataset. We observe that while the linear network shows a period doubling route to chaos, the
 788 ReLU shows long range correlations as observed in real-datasets.

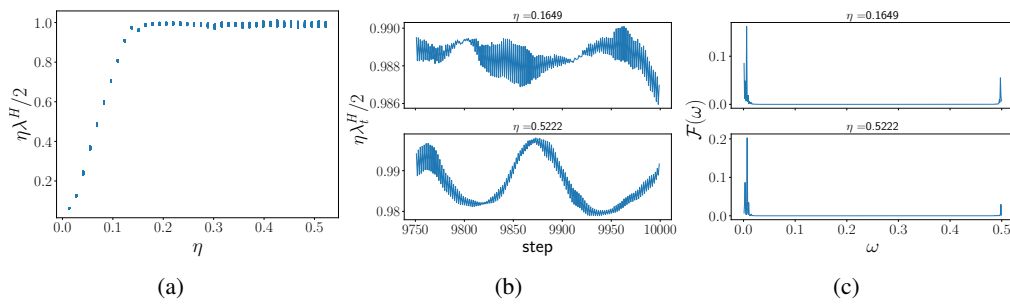


Figure 30: 4-layer ReLU FCNs in SP with $\sigma_w^2 = 0.5$ trained on the generative CIFAR-10 task with MSE loss using GD. Both power spectrums are computed using the last 1000 steps of the corresponding trajectories.