

---

# Mean Field Residual Networks: On the Edge of Chaos

---

**Greg Yang\***  
 Microsoft Research AI  
 gregyang@microsoft.com

**Samuel S. Schoenholz**  
 Google Brain  
 schsam@google.com

## Abstract

We study randomly initialized residual networks using mean field theory and the theory of difference equations. Classical feedforward neural networks, such as those with tanh activations, exhibit exponential behavior on the average when propagating inputs forward or gradients backward. The exponential forward dynamics causes rapid collapsing of the input space geometry, while the exponential backward dynamics causes drastic vanishing or exploding gradients. We show, in contrast, that by adding skip connections, the network will, depending on the nonlinearity, adopt subexponential forward and backward dynamics, and in many cases in fact polynomial. The exponents of these polynomials are obtained through analytic methods and proved and verified empirically to be correct. In terms of the “edge of chaos” hypothesis, these subexponential and polynomial laws allow residual networks to “hover over the boundary between stability and chaos,” thus preserving the geometry of the input space and the gradient information flow. In our experiments, for each activation function we study here, we initialize residual networks with different hyperparameters and train them on MNIST. Remarkably, our *initialization time* theory can accurately predict *test time* performance of these networks, by tracking either the expected amount of gradient explosion or the expected squared distance between the images of two input vectors. Importantly, we show, theoretically as well as empirically, that common initializations such as the Xavier or the He schemes are not optimal for residual networks, because *the optimal initialization variances depend on the depth*. Finally, we have made mathematical contributions by deriving several new identities for the kernels of powers of ReLU functions by relating them to the zeroth Bessel function of the second kind.

## 1 Introduction

Previous works [9, 3, 11] have shown that randomly initialized neural networks exhibit a spectrum of behavior with depth, from stable to chaotic, which depends on the variance of the initializations: the cosine distance of two input vectors converges exponentially fast with depth to a fixed point in  $[0, 1]$ ; if this fixed point is 1, then the behavior is stable; if this fixed point is 0, then the behavior is chaotic. It has been argued in many prior works [1, 9] that effective computation can only be supported by a dynamical behavior that is on **the edge of chaos**. Too much stability prevents the neural network from telling apart two different inputs. While some chaotic behavior can increase the expressivity of a network, too much chaos makes the neural network think two similar inputs are very different. At the same time, the same initialization variances also control how far gradient information can be propagated through the network; the networks with chaotic forward dynamics will tend to suffer from exploding gradients, while networks with stable forward dynamics will tend to suffer from vanishing gradients.

---

\*Work done while at Harvard University

These works have focused on vanilla (fully connected) feedforward networks. Here we consider residual networks [6, 7] (with fully-connected layers and without batchnorm), which are a family of recently proposed neural network architectures that has achieved state-of-the-art performance on image recognition tasks, beating all other approaches by a large margin. The main innovation of this family of architectures is the addition of a passthrough (identity) connection from the previous layer to the next, such that the usual nonlinearity computes the “residual” between the next-layer activation and the previous-layer activation.

In this work, we seek to characterize randomly initialized residual networks. One of our main results is that random residual networks for many nonlinearities such as tanh **live on the edge of chaos**, in that the cosine distance of two input vectors will converge to a fixed point at a polynomial rate, rather than an exponential rate, as with vanilla tanh networks. Thus a typical residual network will slowly cross the stable-chaotic boundary with depth, hovering around this boundary for many layers. In addition, for most of the nonlinearities considered here, the mean field estimate of the gradient grows subexponentially with depth. In fact, for  $\alpha$ -ReLU, the  $\alpha$ th-power of ReLU, for  $\alpha < 1$ , the gradient grows only polynomially. These theoretical results provide some theoretical justification for why residual networks work so well in practice. In our experiments, we are also able to predict surprisingly well the relative performances of *trained* residual networks based only on their initialization hyperparameters, in a variety of settings. In particular, we find that the quality of initialization for tanh resnets is determined by *trainability* (how much gradient explosion on average) while that for ( $\alpha$ -)ReLU resnets is determined by expressivity (how far can two different input vectors be pulled apart) (see Section 6). To the best of our knowledge, this is the first time that a quantity other than gradient explosion/vanishing has been found to control the quality of initialization. We establish theoretically and empirically that the best initialization variances for residual networks depend on the depth of the network (contrary to the feedforward case [11]), so that **common initialization schemes like Xavier [4] or He [5] cannot be optimal**. In fact, even the rationale of He initialization is incorrect for ReLU residual networks because it tries to control gradient dynamics rather than expressivity. However we want to emphasize that we study a simplified model of residual networks in this work, with no batchnorm or convolutional layers, so that these results are not necessarily indicative of the MSRA residual network used in practice [6].

In the body of this paper, we give account of general intuition and/or proof strategy when appropriate for our theoretical results, but we relegate all formal statements and proofs to the appendix.

## 2 Background

Consider a vanilla feedforward neural network of  $L$  layers, with each layer  $l$  having  $N^{(l)}$  neurons; here layer 0 is the input layer. For the ease of presentation we assume all hidden layer widths are the same  $N^{(l)} = N$  for all  $l > 0$ . Let  $x^{(0)} = (x_1^{(0)}, \dots, x_{N^{(0)}}^{(0)})$  be the input vector to the network, and let  $x^{(l)}$  for  $l > 0$  be the activation of layer  $l$ . Then a neural network is given by the equations

$$x_i^{(l)} = \phi(h_i^{(l)}), \quad h_i^{(l)} = \sum_{j=1}^N w_{ij}^{(l)} x_j^{(l-1)} + b_i^{(l)}$$

where (i)  $h^{(l)}$  is the pre-activation at layer  $l$ , (ii)  $w^{(l)}$  is the weight matrix, (iii)  $b^{(l)}$  is the bias vector, and (iv)  $\phi$  is a nonlinearity, for example tanh or ReLU, which is applied coordinatewise to its input.

To lighten up notation, we suppress the explicit layer numbers  $l$  and write

$$x_i = \phi(h_i), \quad h_i = \sum_j w_{ij} \underline{x}_j + b_i$$

where  $\bullet$  implicitly denotes  $\bullet^{(l)}$ , and  $\underline{\bullet}$  denotes  $\bullet^{(l-1)}$  (and analogously,  $\bar{\bullet}$  denotes  $\bullet^{(l+1)}$ ).

A series of papers [9, 10, 11] investigated the “average behavior” of random neural networks sampled via  $w_{ij}^{(l)} \sim \mathcal{N}(0, \sigma_w^2/N)$ ,  $b_i^{(l)} \sim \mathcal{N}(0, \sigma_b^2)$ , for fixed parameters  $\sigma_w$  and  $\sigma_b$ , independent of  $l$ . Consider the expectation of  $\frac{1}{N} \sum_{i=1}^N x_i^2$ , the normalized squared length of  $x$ , over the sampling of  $w$  and  $b$ . Poole et al. [9] showed that this quantity converges to a fixed point exponentially fast for sigmoid nonlinearities. Now suppose we propagate two different vectors  $x^{(0)}$  and  $(x^{(0)})'$  through the

network. Poole et al. [9] also showed that the expectation of the normalized dot product  $\frac{1}{N} \sum_{i=1}^N x_i x'_i$  converges exponentially fast to a fixed point. The ratio between the normalized squared length and the normalized dot product is the cosine distance between  $x$  and  $x'$ . Thus these two exponential convergence results show that the cosine distance converges exponentially fast to a fixed point as well. Intuitively, this means that a vanilla feedforward network “forgets” the geometry of the input space “very quickly,” after only a few layers.

In addition, Schoenholz et al. [11], under certain independence assumptions, showed that the expected normalized squared norm of the gradient also vanishes or explodes in an exponential fashion with depth, with the “half-life” controlled by  $\sigma_w$  and  $\sigma_b$ . They verified that this theoretical “half-life” correlates in practice with the maximal number of layers that are admissible to good performance.

At the same time, Daniely et al. [3] published work of similar nature, but phrased in the language of reproducing kernel Hilbert spaces, and provided high probability estimates that are meaningful for the case when the width  $N$  is finite and the depth is logarithmic in  $N$ . However, they essentially fixed the variance parameters  $\sigma_\bullet$ , and furthermore, their framework (for example the notion of a “skeleton”) does not immediately generalize to the residual network case.

In this work, we show that residual networks have very different dynamics from vanilla feedforward networks. In most cases, the cosine distance convergence rate and the gradient growth rate are subexponential in a residual network, and in most cases, these rates may be polynomial.

### 3 Preliminaries

Residual networks were first introduced by [6] and later refined by [7], and they are now commonplace among deployed neural systems. The key innovation there is the addition of a shortcut connection from the previous layer to the next. We define the following idealized architectures for ease of analysis. Note that we only consider fully-connected affine layers instead of convolutional layers. A **reduced residual network (RRN)** has the recurrence

$$x_i = \phi(h_i) + \underline{x}, \quad h_i = \sum_j w_{ij} \underline{x}_j + b_i.$$

A **(full) residual network (FRN)** in addition has an affine connection given by weights  $v$  and biases  $a$  from the nonlinearity  $\phi(h)$  to the next layer:

$$x_i = \sum_j v_{ij} \phi(h_j) + \underline{x}_i + a_i, \quad h_i = \sum_j w_{ij} \underline{x}_j + b_i$$

We are interested in the “average behavior” of these network when the weights and biases,  $w_{ij}^{(l)}, b_i^{(l)}, v_{ij}^{(l)}$ , and  $a_i^{(l)}$  are sampled i.i.d. from Gaussian distributions resp. with standard deviations  $\sigma_w, \sigma_b, \sigma_v$ , and  $\sigma_a$ , independent from  $l$ . Here we take the variance of  $w_{ij}^{(l)}$  to be  $\sigma_w^2/N$  so that the variance of each  $h_i$  is  $\sigma_w^2$ , assuming each  $\underline{x}_j$  is fixed (similarity for  $v_{ij}^{(l)}$ ). Such an initialization scheme is standard in practice.

We make several key “physical assumptions” to make theoretical computations tractable:

**Axiom 3.1** (Symmetry of activations and gradients). (a) We assume  $\langle (h_i^{(l)})^2 \rangle = \langle (h_j^{(l)})^2 \rangle$  and  $\langle (x_i^{(0)})^2 \rangle = \langle (x_j^{(0)})^2 \rangle$  for any  $i, j, l$ . (b) We also assume that the gradient  $\partial E / \partial x_i^{(l)}$  with respect to the loss function  $E$  satisfies  $\langle (\partial E / \partial x_i^{(l)})^2 \rangle = \langle (\partial E / \partial x_j^{(l)})^2 \rangle$  for any  $i, j, l$ .

One can see that **Axiom 3.1**(a) is satisfied if the input  $x^{(0)} \in \{\pm 1\}^N$  and **Axiom 3.1**(b) is satisfied if **Axiom 3.2** below is true and the gradient at the last layer  $\partial E / \partial x_L \in \{\pm 1\}^N$ . But in general it is justified both empirically and theoretically as an approximation, because  $(h_i^{(l)})^2 - (h_j^{(l)})^2$  stays about constant with  $l$ , but  $(h_i^{(l)})^2$  and  $(h_j^{(l)})^2$  grow rather quickly at the same pace with  $l$  (as will be seen later in calculations), so that their additive difference becomes negligible; similarly for  $(x_i^{(l)})^2$  and  $(\partial E / \partial h_i^{(l)})^2$ .

**Axiom 3.2** (Gradient independence). (a) We assume that we use a different set of weights for back-propagation than those used to compute the network outputs, but sampled i.i.d. from the same distributions. (b) For any loss function  $E$ , we assume that the gradient at layer  $l$ ,  $\partial E/\partial x_i^{(l)}$ , is independent from all activations  $h_j^{(l)}$  and  $x_j^{(l-1)}$  from the previous layer.

Axiom 3.2(a) was first made in [11] for computing the mean field theory of gradients for feedforward tanh networks. This is similar to the practice of feedback alignment [8]. Even though we are the first to explicitly formulate Axiom 3.2(b), in fact it was already applied implicitly in the gradient calculations of [11]. Note that a priori Axiom 3.2(b) is not true, as  $\partial E/\partial x_i^{(l)}$  depends on  $\dot{\phi}(h_k^{(l+1)})$  for every  $k$ , which depend on  $h_j^{(l)}$  for each  $j$ , and which depends on  $x_k^{(l-1)}$  for every  $k$ . Nevertheless, in practice both subassumptions hold very well.

Now we define the central quantities studied in this paper. Inevitably, our paper involves a large amount of notation that may be confusing for the first-time reader. We have included a glossary of symbols (Table A.1) to ameliorate notation confusion.

**Definition 3.3.** Fix an input  $x^{(0)}$ . Define the **length quantities**  $\mathbf{q}^{(l)} := \langle (h_1^{(l)})^2 \rangle$  and  $\mathbf{p}^{(l)} := \langle (x_1^{(l)})^2 \rangle$  for  $l > 0$  and  $\mathbf{p}^{(0)} = \|x^{(0)}\|^2/N$ . Here the expectations  $\langle \bullet \rangle$  are taken over all random initialization of weights and biases for all layers  $l$ , as  $N \rightarrow \infty$  (large width limit).

Note that in our definition, the index 1 does not matter by Axiom 3.1.

**Definition 3.4.** Fix two inputs  $x^{(0)}$  and  $x^{(0)'}$ . We write  $\bullet'$  to denote a quantity  $\bullet$  with respect to the input  $x^{(0)'}$ . Then define the **correlation quantities**  $\gamma^{(l)} := \langle h_1^{(l)} h_1^{(l)'} \rangle$  and  $\lambda^{(l)} := \langle x_1^{(l)} x_1^{(l)'} \rangle$  for  $l > 0$  and  $\gamma^{(0)} = x^{(0)} \cdot x^{(0)'}/N$ , where the expectations  $\langle \bullet \rangle$  are taken over all random initialization of weights and biases for all layers  $l$ , as  $N \rightarrow \infty$  (large width limit). Again, here the index 1 does not matter by Axiom 3.1. By **metric expressivity**, we mean  $\mathbf{s}^{(l)} := \frac{1}{2N} \langle \|x^{(l)} - x^{(l)'}\|^2 \rangle = \frac{1}{2N} (\langle \|x^{(l)}\|^2 \rangle + \langle \|x^{(l)'}\|^2 \rangle - 2\langle x^{(l)} \cdot x^{(l)'} \rangle) = \frac{1}{2} (\mathbf{p}^{(l)} + \mathbf{p}^{(l)'} - \gamma^{(l)})$ . Additionally, define the **cosine distance quantities**  $\mathbf{e}^{(l)} := \gamma^{(l)}/\sqrt{\mathbf{p}^{(l)}\mathbf{p}^{(l)'}}$  and  $\mathbf{c}^{(l)} := \lambda^{(l)}/\sqrt{\mathbf{q}^{(l)}\mathbf{q}^{(l)'}}$ , and we will also call  $\mathbf{e}^{(l)}$  **angular expressivity**.

In this paper, for the ease of presentation, we assume  $\mathbf{p}^{(0)} = \mathbf{p}^{(0)'}$ . Then, as we will see,  $\mathbf{p}^{(l)} = \mathbf{p}^{(l)'}$ ,  $\mathbf{q}^{(l)} = \mathbf{q}^{(l)'}$  for all  $l$ , and as a result,  $\mathbf{e}^{(l)} = \gamma^{(l)}/\mathbf{p}^{(l)}$  and  $\mathbf{s}^{(l)} = \mathbf{p}^{(l)} - \gamma^{(l)} = (1 - \mathbf{e}^{(l)})\mathbf{p}^{(l)}$ .

**Definition 3.5.** Fix an input  $x^{(0)}$  and a gradient vector  $(\partial E/\partial x_i^{(L)})_i$  of some loss function  $E$  with respect to the last layer  $x^{(L)}$ . Then define the **gradient quantities**  $\chi^{(l)} := \langle (\partial E/\partial x_1^{(l)})^2 \rangle$ ,  $\chi_\bullet^{(l)} := \langle (\partial E/\partial \bullet_1^{(l)})^2 \rangle$  for  $\bullet = a, b$ , and  $\chi_\bullet^{(l)} := \langle (\partial E/\partial \bullet_1^{(l)})^2 \rangle$  for  $\bullet = w, v$ . Here the expectations are taken with Axiom 3.2 in mind, over both random initialization of forward and backward weights and biases, as  $N \rightarrow \infty$  (large width limit). Again, the index 1 or 11 does not matter by Axiom 3.1.

**Asymptotic notations.** The expressions  $f = O(g) \iff g = \Omega(f)$  have their typical meanings, and  $f = \Theta(g)$  iff  $f = O(g), g = O(f)$ . We take  $f(x) = \tilde{O}(g(x)) \iff g(x) = \tilde{\Omega}(f(x))$  to mean  $f(x) = O(g \log^k x)$  for some  $k \in \mathbb{Z}$  (this is slightly different from the standard usage of  $\tilde{O}$ ), and  $f = \tilde{\Theta}(g) \iff f = \tilde{O}(g) \ \& \ g = \tilde{O}(f)$ . We introduce a new notation:  $f = \check{\Theta}(g)$  if  $f(x) = O(g(x) \cdot x^\epsilon)$  and  $f(x) = \Omega(g(x) \cdot x^{-\epsilon})$ , as  $x \rightarrow \infty$ , for any  $\epsilon > 0$ . All asymptotic notations are sign-less, i.e. can indicate either positive or negative quantities, unless stated otherwise.

## 4 Overview

The primary reason we may say anything about the average behavior of any of the above quantities is the central limit theorem: every time the activations of the previous layer pass through an affine layer whose weights are sampled i.i.d., the output is a sum of a large number of random variables, and thus follows approximately Gaussian distributions. The mean and variance of these distributions can be computed by keeping track of the mean and variances of the activations in the previous layer.

In what follows, we use this technique to derive recurrence equations governing  $\mathbf{p}, \mathbf{q}, \gamma, \lambda, \chi$  for different architectures and different activation functions. We use these equations to investigate the



dynamics of  $\mathbf{e}$  and  $\mathbf{s}$ , the key quantities in the forward pass, and the dynamics of  $\chi$ , the key quantity in the backward pass.

The cosine distance  $\mathbf{e}$  in some sense measures the angular geometry of two vectors. If  $\mathbf{e} = 1$ , then the vectors are parallel; if  $\mathbf{e} = 0$ , then they are orthogonal. Just as in [9] and [11], we will show that in all of the architectures and activations we consider in this paper,  $\mathbf{e}^{(l)}$  converges to a fixed point  $\mathbf{e}^*$  as  $l \rightarrow \infty$ <sup>1</sup>. Thus, on the average, as vectors propagate through network, the geometry of the original input space, for example, linear separability, is “forgotten” by residual networks as well as by vanilla networks. But we will prove and verify experimentally that, while Poole et al. [9] and [11] showed that the convergence rate to  $\mathbf{e}^*$  is exponential in a vanilla network, the convergence rate is rather only polynomial in residual networks, for tanh and  $\alpha$ -ReLU (Defn 5.2) nonlinearities; see Thm B.5, Thm B.11, Thm B.17, and Thm B.18. This slow convergence preserves geometric information in the input space, and allows a typical residual network to “hover over the edge of chaos”: Even when the cosine distance  $\mathbf{e}^{(l)}$  converges to 0, corresponding to “chaos”, (resp. 1, corresponding to “stability”), for the number of layers usually seen in practice,  $\mathbf{e}^{(l)}$  will reside well away from 0 (resp. 1).

Similarly, the quantity  $\mathbf{s}$  measures the metric geometry of two vectors. The evolution of  $\mathbf{s}^{(l)}$  with  $l$  tells us the ability of the average network to separate two input points in terms of Euclidean distance. Again, for tanh and  $\alpha$ -ReLU ( $\alpha < 1$ ) nonlinearities,  $\mathbf{s}$  varies only polynomially with  $l$ .

On the other hand,  $\chi^{(l)}$  measures the size of gradient at layer  $l$ , and through it we track the dynamics of gradient backpropagation, be it explosion or vanishing. In contrast to vanilla tanh networks, which can experience both of these two phenomenon depending on the initialization variances, typical residual networks cannot have vanishing gradient, in the sense of vanishing  $\chi^{(l)}$  as  $l \rightarrow 1$ ; see Thm B.5 and Thm B.12. Furthermore, while vanilla tanh networks exhibit exponentially vanishing or exploding gradients, all of the activation/architecture pairings considered here, except the full residual network with ReLU, have subexponential gradient dynamics. While tanh residual networks (reduced or full) has  $\chi^{(0)} \approx \exp(\Theta(\sqrt{l}))\chi^{(l)}$  (Thm B.13),  $\alpha$ -ReLU residual networks for  $\alpha < 1$  have  $\chi^{(0)} \approx \text{poly}(l)\chi^{(l)}$  (Thm B.20). Instead of  $\partial E/\partial x_i$ , we may also consider the size of gradients of actual trainable parameters. For tanh and  $\alpha$ -ReLU with  $\alpha < 1$ , they are still subexponential and polynomial (Thm B.21). On the other hand, while  $\chi^{(0)} = \exp(\Theta(l))\chi^{(l)}$  for a ReLU resnet, its weight gradients have size independent of layer, within  $O(1)$  (Thm B.21)! This is the only instance in this paper of gradient norm being completely preserved across layers.

The above overviews the theoretical portion of this paper. Through experiments, we discover that we can very accurately predict whether one random initialization leads to better performance than another on the test set, after training, by leveraging this theory we build. Residual networks of different nonlinearities have different *controlling quantities*: for resnets with tanh, the optimal initialization is obtained by controlling the gradient explosion  $\chi^{(0)}/\chi^{(L)}$ ; whereas for ReLU and  $\alpha$ -ReLU, the optimal initialization is obtained by maximizing  $\mathbf{s}$  without running into numerical issues (with floating point computation). See Section 6 for details.

Over the course of our investigation of  $\alpha$ -ReLU, we derived several new identities involving the associated kernel functions, first defined in [2], which relate them to the zeroth Bessel functions (Lemmas C.31 to C.34).

## 5 Theoretical Results

In what follows in the main text, we assume  $\sigma_{\bullet} > 0$  for all  $\bullet = w, v, b, a$ ; in the appendix, the formal statement of each main theorem will contain results for other cases. We are interested in the two major categories of nonlinearities used today: tanh-like and rectified units. We make the following formal definitions as a foundation for further consideration.

**Definition 5.1.** We say a function  $\phi$  is **tanh-like** if  $\phi$  is antisymmetric ( $\phi(-x) = -\phi(x)$ ),  $|\phi(x)| \leq 1$  for all  $x$ ,  $\phi(x) \geq 0, \forall x \geq 0$ , and  $\phi(x)$  monotonically increases to 1 as  $x \rightarrow \infty$ .

**Definition 5.2.** Define the  $\alpha$ -ReLU  $\psi_{\alpha}(x) = x^{\alpha}$  if  $x > 0$  and 0 otherwise.<sup>2</sup>

By applying the central limit theorem as described in the last section, we derive a set of recurrences for different activation/architecture pairs, shown in Table 1 (see appendix for proofs). They leverage certain integral transforms<sup>3</sup> as in the following

**Table 1:** Main Recurrences

Antisymmetric/RRN		Any/FRN	
$\mathbf{q} = \sigma_w^2 \underline{\mathbf{p}} + \sigma_b^2$	$\mathbf{p} = V\phi(\mathbf{q}) + \underline{\mathbf{p}}$	$\mathbf{q} = \sigma_w^2 \underline{\mathbf{p}} + \sigma_b^2$	$\mathbf{p} = \sigma_v^2 V\phi(\mathbf{q}) + \sigma_a^2 + \underline{\mathbf{p}}$
$\lambda = \sigma_w^2 \underline{\gamma} + \sigma_b^2$	$\gamma = W\phi(\mathbf{q}, \lambda) + \underline{\gamma}$	$\lambda = \sigma_w^2 \underline{\gamma} + \sigma_b^2$	$\gamma = \sigma_v^2 W\phi(\mathbf{q}, \lambda) + \sigma_a^2 + \underline{\gamma}$
	$\chi = (\sigma_w^2 V\dot{\phi}(\mathbf{q}) + 1)\chi$		$\chi = (\sigma_v^2 \sigma_w^2 V\dot{\phi}(\mathbf{q}) + 1)\chi$
Theorems <a href="#">B.2</a> , <a href="#">B.3</a> , <a href="#">B.5</a>		Theorems <a href="#">B.8</a> , <a href="#">B.10</a> , <a href="#">B.12</a>	

**Table 2:** Summary of Main Dynamics Results. Note that while  $\chi^{(l)}$  is exponential for ReLU/FRN, the gradients with respect to weight parameters have norms ( $\chi_w$  and  $\chi_v$ ) constant in  $l$  ([Thm B.21](#)). Also, the  $\chi^{(l)}$  entry for  $\alpha$ -ReLU is for  $\alpha \in (3/4, 1)$  only

	Tanh/RRN	Tanh/FRN	ReLU/FRN	$\alpha$ -ReLU/FRN, $\alpha < 1$
$\mathbf{p}^{(l)}$	$\Theta(l)$ , <a href="#">B.2</a>	$\Theta(l)$ , <a href="#">B.9</a>	$\exp(\Theta(l))$ , <a href="#">B.16</a>	$\Theta(l^{1/(1-\alpha)})$ , <a href="#">B.16</a>
$\mathbf{s}^{(l)}$	$\Theta(l)$ , <a href="#">B.4</a>	$\Theta(l)$ , <a href="#">B.11</a>	$\exp(\Theta(l))$ , <a href="#">B.17</a>	$\Theta(l^{1/(1-\alpha)})$ , <a href="#">B.18</a>
$\mathbf{e}^{(l)} - \mathbf{e}^*$	$\check{\Theta}(l^{\frac{2}{\pi}-1})$ , <a href="#">B.4</a>	poly( $l$ ), <a href="#">B.11</a>	$\Theta(l^{-2})$ , <a href="#">B.17</a>	poly( $l$ ), <a href="#">B.18</a>
$\chi^{(l)}$	$\exp(\Theta(\sqrt{l}))$ , <a href="#">B.6</a>	$\exp(\Theta(\sqrt{l}))$ , <a href="#">B.12</a>	$\exp(\Theta(l))$ , <a href="#">B.20</a>	$\Theta(l^{\frac{\alpha^2}{(1-\alpha)(2\alpha-1)}})$ , <a href="#">B.20</a>

**Definition 5.3.** Define the transforms  $V$  and  $W$  by  $V\phi(q) := \mathbb{E}[\phi(z)^2 : z \sim \mathcal{N}(0, q)]$  and  $W\phi(\rho, \nu) := \mathbb{E}[\phi(z)\phi(z') : (z, z') \sim \mathcal{N}(0, \begin{pmatrix} \rho & \nu \\ \nu & \rho \end{pmatrix})]$ .

These recurrences are able to track the corresponding quantities in practice very well. For example, [Fig. 1](#) compares theory vs experiments for the tanh/FRN pair. The agreement is very good for tanh/RRN (not shown, but similar to the case of tanh/FRN with  $\sigma_v = 1$  and  $\sigma_a = 0$ ) and  $\alpha$ -ReLU/FRN as well (see [Fig. A.1](#)).

As mentioned in previous sections, we seek to characterize the long term/high depth behavior of all of the quantities defined in [Section 2](#). To do so, we solve for the asymptotics of the recurrences in [Table 1](#), where  $\phi$  is instantiated with tanh or  $\alpha$ -ReLU. Our main dynamics results are summarized in [Table 2](#).

## 5.1 Tanh

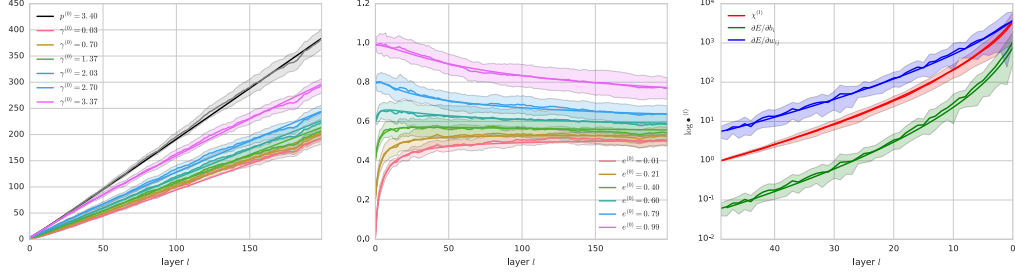
**Forward dynamics.** When  $\phi = \tanh$ ,  $\mathbf{p}^{(l)}$  and  $\mathbf{q}^{(l)}$  increase as  $\Theta(l)$  in either RRN or FRN ([Thm B.2](#)), as one might expect by observing that  $V \tanh(\mathbf{q}) \rightarrow 1$  as  $\mathbf{q} \rightarrow \infty$  so that, for example in the RRN case, the recurrence  $\mathbf{p} = V \tanh(\mathbf{q}) + \underline{\mathbf{p}}$  becomes  $\mathbf{p} = 1 + \underline{\mathbf{p}}$ . This is confirmed graphically by the black lines of the leftmost chart of [Fig. 1](#). We carefully verify that this intuition is correct in its proof in the appendix, and find that in fact  $\mathbf{p}^{(l)} \sim l$  in the RRN case and  $\mathbf{p}^{(l)} \sim (\sigma_v^2 + \sigma_a^2)l$  in the FRN case.

What about  $\gamma^{(l)}$ ? The middle chart of [Fig. 1](#) shows that over time,  $\mathbf{e}^{(l)} = \gamma^{(l)}/\mathbf{p}^{(l)}$  contracts toward the center of the interval  $[0, 1]$ , but from the looks of it, it is not clear whether there is a stable fixed point  $\mathbf{e}^*$  of  $\mathbf{e}$  or not. We prove that, in fact, **all trajectories of  $\mathbf{e}$  not starting at 1 do converge to a single fixed point, but only at a polynomial rate**, in both the RRN and FRN cases ([Thm B.2](#) and [Thm B.10](#)); we can even explicitly compute the fixed point and the rate of convergence: For FRN, there is a **unique stable fixed point**  $\mathbf{e}^* < 1$  determined by the equation

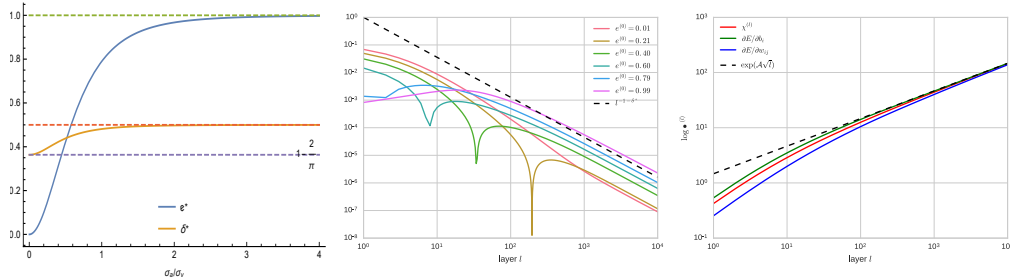
$$\mathbf{e}^* = \frac{1}{\sigma_v^2 + \sigma_a^2} [\sigma_v^2 \frac{2}{\pi} \arcsin(\mathbf{e}^*) + \sigma_a^2],$$

and  $|\mathbf{e}^{(l)} - \mathbf{e}^*|$  decreases like  $l^{-\delta^*}$ , where

$$\delta^* := 1 - \frac{2}{\pi} \frac{1}{\sqrt{1 - (\mathbf{e}^*)^2}} \frac{\sigma_v^2}{\sigma_v^2 + \sigma_a^2}.$$



**Figure 1:** Our equations predict the relevant quantities very well in practice. These plots make the comparison between prediction and measurements for the full resnet with tanh activation, with  $\sigma_v^2 = 1.5$ ,  $\sigma_a^2 = .5$ ,  $\sigma_w^2 = 1.69$ ,  $\sigma_b^2 = .49$ . Left-to-right: **(a)**  $\mathbf{p}^{(l)}$  and  $\gamma^{(l)}$  against layer  $l$  for 200 layers. **(b)**  $\mathbf{e}^{(l)} = \gamma^{(l)}/\mathbf{p}^{(l)}$  against  $l$  for 200 layers. Both (a) and (b) trace out curves for different initial conditions. **(c)** Different gradient quantities against  $l$  for 50 layers. From left to right the layer number  $l$  decreases, following the direction of backpropagation. Notice that the gradient increases in norm as  $l \rightarrow 1$ . All three figures exhibit smooth curves, which are theoretical estimates, and irregular curves with shades around them, which indicate empirical means and standard deviations (both of which taken in regular scale, not log scale). (a) and (b) are made with 20 runs of resnets of width 1000. (c) is made with 25 runs of resnets of width 250.



**Figure 2:** Left-to-right: **(a)** Plots of  $\mathbf{e}^*$  and  $\delta^*$  against  $\sigma_a/\sigma_v$ . **(b)** In log-log scale: the dashed line is  $l^{-\delta^*-1}$ , and the colored lines are  $\mathbf{e}^{(l)} - \mathbf{e}^{(l-1)}$  for different initial conditions  $\mathbf{e}^{(0)}$ . That they become parallel at about  $l = 400$  on verifies that  $\mathbf{e}^{(l)} = \Theta(l^{-\delta^*})$ .<sup>4</sup> **(c)** In log-log scale: The dashed line is  $\mathcal{A}\sqrt{l}$  ( $\mathcal{A}$  given in [Thm B.13](#)), and the colored lines are  $\log(\bullet^{(1)}/\bullet^{(l)})$  for  $\bullet = \chi, \chi_b, \chi_w$ . That they all converge together starting around  $l = 1000$  indicates that the approximation in [Thm B.13](#) is very good for large  $l$ .

Since  $\mathbf{e}^* < 1$ ,  $\mathbf{s} = (1 - \mathbf{e})\mathbf{p} = \Theta(\mathbf{p}) = \Theta(l)$ . The case of RRN can be viewed as a special case of the above, setting  $\sigma_v^2 = 1$  and  $\sigma_a^2 = 0$ , which yields  $\mathbf{e}^* = 0$  and  $\delta^* = 1 - \frac{2}{\pi}$ . We observe that both  $\mathbf{e}^*$  and  $\delta^*$  only depend on the ratio  $\rho := \sigma_a/\sigma_v$ , so in [Fig. 2](#) we graph these two quantities as a function of  $\rho$ .  $\mathbf{e}^*$  and  $\delta^*$  both increase with  $\rho$  and asymptotically approach 1 and  $1/2$  respectively from below. When  $\rho = \sigma_a = 0$ ,  $\mathbf{e}^* = 0$  and  $\delta^* = 1 - \frac{2}{\pi}$ . Thus the rate of convergence at its **slowest** for tanh/FRN is  $\delta^* = 1 - \frac{2}{\pi} \approx 0.36338$ , where asymptotically the network tends toward a **chaotic regime**  $\mathbf{e}^* = 0$ , corresponding to a large weight variance and a small bias variance; it at its **fastest** is  $\delta^* = 1/2$ , where asymptotically the network tends toward a **stable regime**  $\mathbf{e}^* = 1$ , corresponding to a large bias variance and small weight variance. We verify  $\delta^*$  by comparing  $\mathbf{e}^{(l)} - \mathbf{e}^{(l-1)}$  to  $l^{-\delta^*-1}$  in log-log scale. If  $\mathbf{e}^{(l)} = \Theta(l^{-\delta^*})$ , then  $\mathbf{e}^{(l)} - \mathbf{e}^{(l-1)} = \Theta(l^{-\delta^*-1})$  and should obtain the same slope as  $l^{-\delta^*-1}$  as  $l \rightarrow \infty$ . The middle figure of [Fig. 2](#) ascertains that this is indeed the case, starting around layer number 400.

**Backward dynamics.** Finally, we show that the gradient is approximated by

$$\chi^{(m)} = \exp(\mathcal{A}(\sqrt{l} - \sqrt{m}) + O(\log l - \log m))\chi^{(l)} \quad (\star)$$

where  $\mathcal{A} = \frac{4}{3}\sqrt{\frac{2}{\pi}}\sigma_w$  in the RRN case and  $\mathcal{A} = \frac{4}{3}\sqrt{\frac{2}{\pi}}\frac{\sigma_w^2\sigma_w}{\sqrt{\sigma_v^2 + \sigma_a^2}}$  in the FRN case ([Thm B.6](#) and [Thm B.13](#)). The rightmost plot of [Fig. 2](#) verifies that indeed, for large  $l \geq 1000$ , this is a very good approximation. This demonstrates that the mean field assumption of independent backpropagation weights is very practical and convenient even for residual networks.

Note that in the FRN case, the constant  $\mathcal{A}$  can be decomposed into  $\mathcal{A} = \frac{4}{3}\sqrt{\frac{2}{\pi}} \cdot \sigma_v \cdot \sigma_w \cdot (1 + \sigma_a^2/\sigma_v^2)^{-1/2}$ . Consider the ratio  $\rho := \sigma_a/\sigma_v$ . If  $\rho \gg 1$ , then  $\mathbf{e}^* \approx 1$  (Fig. C.17), meaning that the typical network essentially computes a constant function, and thus unexpressive; at the same time, large  $\rho$  makes  $\mathcal{A}$  small, and thus ameliorating the gradient explosion problem, making the network more trainable. On the other hand, if  $\rho \ll 1$ , then  $\mathbf{e}^* \approx 0$  (Fig. C.17), the typical network can tease out the finest differences between any two input vectors, and a final linear layer on top of such a network should be able to express a wide variety of functions [9]; at the same time, small  $\rho$  increases  $\mathcal{A}$ , worsening the gradient explosion problem, making the network less trainable. This is the same expressivity-trainability tradeoff discussed in [11].

## 5.2 $\alpha$ -ReLU

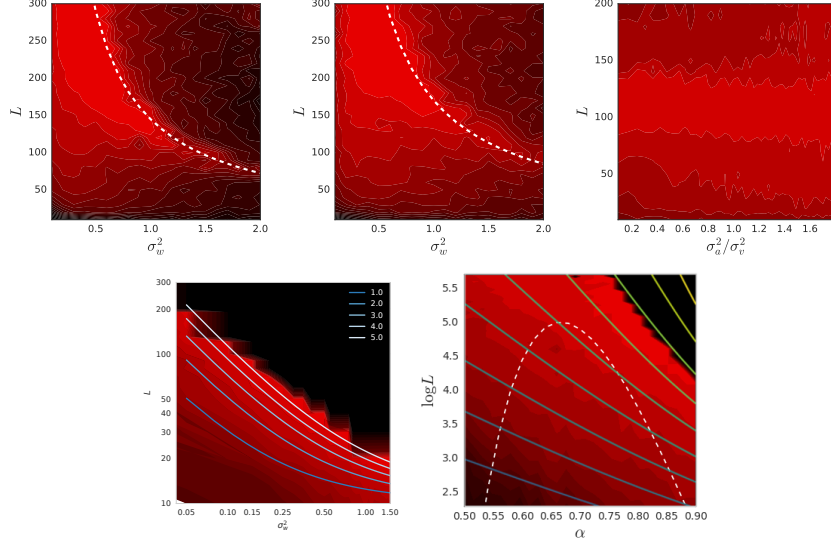
**Forward dynamics.** As with the tanh case, to deduce the asymptotic behavior of random  $\alpha$ -ReLU resnets, we need to understand the transforms  $V\psi_\alpha$  and  $W\psi_\alpha$ . Fortunately,  $V\psi_\alpha$  has a closed form, and  $W\psi_\alpha$  has been studied before [2]. In particular, if  $\alpha > -\frac{1}{2}$ , then  $V\psi_\alpha(\mathbf{q}) = c_\alpha \mathbf{q}^\alpha$ , where  $c_\alpha$  is a constant with a closed form given by Lemma B.15. In addition, by [2], we know that  $W\psi_\alpha(\mathbf{q}, \mathbf{c}\mathbf{q}) = V\psi_\alpha(\mathbf{q})\mathbb{J}_\alpha(\mathbf{c})$  for  $\mathbb{J}_\alpha$  given in Appendix C.7.1. Fig. C.17 shows a comparison of  $\mathbb{J}_\alpha$  for different  $\alpha$ s along with the identity function.

Substituting in  $c_\alpha \mathbf{q}^\alpha$  for  $V\psi_\alpha$ , we get a difference equation  $\mathbf{p} - \underline{\mathbf{p}} = \sigma_v^2 c_\alpha (\sigma_w^2 \underline{\mathbf{p}} + \sigma_b^2)^\alpha + \sigma_a^2$  governing the evolution of  $\mathbf{p}$ . This should be reminiscent of the differential equation  $\dot{P}(l) = CP(l)^\alpha$ , which has solution  $\propto l^{1/(1-\alpha)}$  for  $\alpha < 1$ , and  $\propto \exp(Cl)$  when  $\alpha = 1$ . And indeed, the solutions  $\mathbf{p}^{(l)}$  to these difference equations behave asymptotically exactly like so (Thm B.16). Thus **ReLU behaves very explosively compared to  $\alpha$ -ReLU with  $\alpha < 1$** . In fact, in simulations, for  $\sigma_w^2 = 1.69$  and  $\sigma_v^2 = 1.5$ , the ReLU resnets overflows into infs after around 100 layers, while there's no problem from any other kind of networks we consider.

Regardless,  **$\alpha$ -ReLU for all  $\alpha$  messages  $\mathbf{e}^{(l)}$  toward a fixed point  $\mathbf{e}^*$  that depends on  $\alpha$** . When  $\phi = \psi_1$ , the standard ReLU,  $\mathbf{e}^{(l)}$  converges to 1 asymptotically as  $Cl^{-2}$  for an explicit constant  $C$  depending on  $\sigma_v$  and  $\sigma_w$  only (Thm B.17), so that  $\mathbf{s} = (1 - \mathbf{e})\mathbf{p} = \Theta(l^{-2} \exp(\Theta(l))) = \exp(\Theta(l))$ . When  $\phi = \psi_\alpha$  for  $\alpha < 1$ , then  $\mathbf{e}^{(l)}$  converges to the nonunit fixed point  $\mathbf{e}^*$  of  $\mathbb{J}_\alpha$  at a rate of  $\Theta(l^{-\mu})$ , where  $\mu = (1 - \mathbb{J}_\alpha(\mathbf{e}^*)) / (1 - \alpha)$  is independent of the variances (Thm B.18), so that  $\mathbf{s} = \Theta(\mathbf{p})$ . These rates are verified in Fig. A.2.

**Backward dynamics.** Finally, we have also characterized the rate of gradient growth for any  $\alpha \in (\frac{3}{4}, 1]$ .<sup>5</sup> **In the case of  $\alpha = 1$ , the dynamics of  $\chi$  is exponential**, the same as that of  $\mathbf{p}$ ,  $\chi^{(l-m)} = \chi^{(l)} B^m$  where  $B = \frac{1}{2}\sigma_v^2\sigma_w^2 + 1$ . **For  $\alpha \in (\frac{3}{4}, 1)$ , the dynamics is polynomial**, but with different exponent in general from that of the forward pass:  $\chi^{(l-m)} = \Theta(1)\chi^{(l)}(l/(l-m))^R$  for  $R = \frac{\alpha^2}{(1-\alpha)(2\alpha-1)}$ , where the constants in  $\Theta(1)$  do not depend on  $l$  or  $m$ . This exponent  $R$  is minimized on  $\alpha \in [\frac{3}{4}, 1)$  at  $\alpha = 3/4$ , where  $R = 9/2$  (but on  $\alpha \in (\frac{1}{2}, 1)$  it is minimized at  $\alpha = 2/3$ , where  $R = 4$ ); see Fig. B.8. These exponents are verified empirically in Fig. A.2.

Looking only at  $\chi$  and the gradients against the biases, it seems that ReLU suffers from a dramatic case of exploding gradients. But in fact, because  $\chi$  gains a factor of  $B$  moving backwards while  $\mathbf{p}$  loses a factor of  $B$ , the gradient norm  $\chi_w^{(l-m)}$  (and similarly for  $\chi_v^{(l-m)}$ ) is independent of how far,  $m$ , the gradient has been propagated (Thm B.21) — this is certainly the best gradient preservation among all of the models considered in this paper. Thus strangely, random ReLU FRN exhibits both the best (constant for  $v$  and  $w$ ) and the worse (exponential for  $a$  and  $b$ ) gradient dynamics. This begs the question, then, is this a better deal than other  $\alpha$ -ReLU for which for any learnable parameter we have at most a polynomial blowup with depth in its gradient? Our experiments (discussed below) show that  $\alpha$ -ReLU is useful to the extent that smaller  $\alpha$  avoids numerical issues with exponentiating forward and backward dynamics, but the best performance is given by the largest  $\alpha$  that avoids them (Fig. 3(c, d)); in fact, the metric expressivity  $\mathbf{s}$ , determines performance, not gradient explosion (see  $\alpha$ -ReLU experiments).



**Figure 3:** From left to right, top to bottom: (a) and (b):  $\sigma_w^2$ ,  $L$ , and test set accuracy of a grid of tanh reduced (left) and full (right) resnets trained on MNIST. Color indicates performance, with higher colors indicating higher accuracy on test set. Other than the values on the axes, we have fixed  $\sigma_b^2 = \sigma_a^2 = \frac{1}{2}$  and  $\sigma_v^2 = 1$ . The white dotted lines are given by  $\sigma_w^2 L = C$ , where  $C = 170$  on the left and  $C = 145$  on the right. We see that both dotted lines accurately predict the largest optimal  $\sigma_w$  for each depth  $L$ . (c) Varying the ratio  $\sigma_a^2/\sigma_v^2$  while fixing  $\sigma_v/\sqrt{1 + \sigma_a^2/\sigma_v^2}$ , and thus fixing  $\mathcal{A}$ , the leading constant of  $\log \chi^{(0)}/\chi^{(L)}$ . (d) in log-log scale: Heatmap gives the test accuracies of ReLU FRN for varying  $\sigma_w^2$  and  $L$ . Curves give level sets for the log ratios  $\log \mathbf{s}^{(L)}/\mathbf{s}^{(0)} \approx \log \mathbf{p}^{(L)}/\mathbf{p}^{(0)} \approx \log \chi^{(0)}/\chi^{(L)} = L \log(1 + \sigma_v^2 \sigma_w^2/2)$ . (e) Red heatmap shows the test accuracies of a grid of  $\alpha$ -ReLU FRN with varying  $\alpha$  and  $L$  as shown, but with all  $\sigma_\bullet$ s fixed. The white dashed curve gives a typical contour line of  $L^R = \text{const}$ , where  $R = \frac{\alpha^2}{(1-\alpha)(2\alpha-1)}$ . The yellow-to-blue curves form a set of level curves for  $\mathbf{s}^{(l)} = \mathbf{p}^{(l)} - \gamma^{(l)} = \text{const}$ , with yellow curves corresponding to higher levels.

## 6 Experimental Results

Our experiments show a dichotomy of what matters in initialization: for tanh resnets, quality of an initialization is determined by how much gradient explosion there is (measured by  $\chi^{(0)}/\chi^{(L)}$ ); for ( $\alpha$ )-ReLU resnets, it is determined by how expressive the random network is (measured by the metric expressivity  $\mathbf{s}^{(L)}$ ). We hypothesize this is because in tanh resnets, the gradient dynamics is much more explosive than the expressivity dynamics ( $\exp(\Theta(\sqrt{l}))$  vs  $\Theta(l)$ ), whereas for ReLU it's somewhat the opposite ( $\chi_w, \chi_v = \Theta(1)$  vs  $\mathbf{s} = \exp(\Theta(l))$ ).

**Tanh, vary  $\sigma_w$ .** We train a grid of reduced and full tanh resnets on MNIST, varying the variance  $\sigma_w^2$  and the number of layers (for FRN we fix  $\sigma_v = 1$ ). The results are indicated in Fig. 3(a, b). We see that in either model, deeper resnets favor much smaller  $\sigma_w$  than shallower ones. The white dotted lines in Fig. 3(a, b) confirm our theory: according to Eq. (\*), for the same gradient ratio  $R = \chi^{(0)}/\chi^{(L)}$ , we want  $\log R \approx \sigma_w \sqrt{L}$ . Indeed, the white dotted lines in Fig. 3(a, b) trace out such a level curve and it remarkably pinpoints the largest  $\sigma_w$  that gives the optimal test set accuracy for each depth  $L$ . Why isn't the best initialization given by  $R = 1 \iff \sigma_w = 0$ ? We believe that when  $L$  and/or  $\sigma_w$  is small, gradient dynamics no longer dominates the initialization quality because it has "less room to explode," and expressivity issues start to dampen the test time performance.

**Tanh, vary  $\sigma_a^2/\sigma_v^2$ .** As suggested in the analysis of Eq. (\*), the ratio  $\rho^2 = \sigma_a^2/\sigma_v^2$  determines the fixed point  $\mathbf{e}^*$  and its convergence rate by itself while also contributes to the rate of gradient explosion in tanh FRN. We seek to isolate its effect on forward dynamics by varying  $\sigma_v$  with  $\rho$  such that  $\sigma_v/\sqrt{1 + \rho^2}$  is kept constant, so that the leading term of the log gradient ratio is kept approximately equal for each  $L$  and  $\rho$ . Fig. 3(c) shows the test accuracies of a grid of tanh FRN initialized with such an ensemble of  $\sigma_\bullet$ s. What stands out the most is that performance is maximized essentially

around a fixed value of  $L$  regardless of  $\rho$ , which shows that indeed gradient dynamics determines the initialization quality in tanh resnets. There is also a minor increase in performance with increasing  $\rho$  regardless of  $L$ ; this is counterintuitive as increasing  $\rho$  means “decreasing expressivity.” It is currently not clear what accounts for this effect.

**ReLU, vary  $\sigma_w$ .** We train a grid of ReLU FRN on MNIST, varying  $\sigma_w^2 \in [0, 1.5]$  while fixing  $\sigma_v^2 = 1, \sigma_a^2 = \sigma_b^2 = \frac{1}{2}$ . The resulting test set accuracies are shown in Fig. 3(d). The dark upper region signifies failure of training caused by numerical issues with exploding activation and gradient norms: This corresponds to the region where  $\mathbf{p}^{(L)}$ , which is a measure of the mean magnitude of an neuronal activation in layer  $L$ , becomes too big. We see that the best test accuracies are given by depths just below where these numerical issues occur. However, if we were to predict that the optimal init is the one minimizing  $\chi^{(0)}/\chi^{(L)} \geq 1$ , then we would be wrong — in fact it is exactly the opposite. In this case, the dynamics of  $\mathbf{s}^{(l)}, \mathbf{p}^{(l)}$ , and  $\chi^{(0)}/\chi^{(l)}$  are approximately the same (all  $\exp(\Theta(l))$  with the same hidden constants), and optimal performance corresponds to the highest  $\mathbf{s}^{(L)}, \mathbf{p}^{(L)}$ , and  $\chi^{(0)}/\chi^{(L)}$  without running into infs.

**$\alpha$ -ReLU, vary  $\alpha$ .** We similarly trained a grid of  $\alpha$ -ReLU FRN on MNIST, varying only  $\alpha$  and the depth, fixing all  $\sigma$ . Fig. 3(e) shows their test accuracies. We see similar behavior to ReLU, where when the net is too deep, numerical issues doom the training (black upper right corner), but the best performance is given by  $L$  just below where this problem occurs. In this case, if we were to predict optimality based on minimizing gradient explosion, we would be again wrong, and furthermore, the contour plot of  $\chi^{(0)}/\chi^{(L)}$  (white dashed line) now gives no information at all on the test set accuracy. In contrast, the contours for  $\mathbf{s}^{(l)}$  succeeds remarkably well at this prediction (yellow/green lines).<sup>6</sup> By interpolation, this suggests that indeed in the ReLU case, it is expressivity, not trainability, which determines performance at test time.

In all of our experiments, we did not find e dynamics to be predictive of neural network performance.

## 7 Conclusion

In this paper, we have extended the mean field formalism developed by [9, 10, 11] to residual networks, a class of models closer to practice than classical feedforward neural networks as were investigated earlier. We proved and verified that in both the forward and backward passes, most of the residual networks discussed here do not collapse their input space geometry or the gradient information exponentially. We found our theory incredibly predictive of test time performance despite saying nothing about the dynamics of training. In addition, we overwhelmingly find, through theory and experiments, that an optimal initialization scheme must take into account the depth of the residual network. The reason that Xavier [4] or He [5] scheme are not the best for residual networks is in fact not that their statistical assumptions are fragile — theirs are similar to our mean field theoretic assumptions, and they hold up in experiments for large width — but rather that their structural assumptions on the network break very badly on residual nets.

**Open Problems.** Our work thus have shown that optimality of initialization schemes can be very unstable with respect to architecture. We hope this work will form a foundation toward a mathematically grounded initialization scheme for state-of-the-art architectures like the original He et al. residual network. To do so, there are still two major components left to study out of the following three: 1. Residual/skip connection 2. Batchnorm 3. Convolutional layers. Recurrent architectures and attention mechanisms are also still mostly unexplored in terms of mean field theory. Furthermore, many theoretical questions still yet to be resolved; the most important with regard to mean field theory is: why can we make Axioms 3.1 and 3.2 and still be able to make accurate predictions? We hope to make progress on these problems in the future and encourage readers to take part in this effort.



## Acknowledgments

Thanks to Jeffrey Ling for early exploration experiments and help with the initial draft. Thanks to Felix Wong for offering his wisdom and experience working in statistical physics.

## References

- [1] Nils Bertschinger and Thomas Natschlger. Real-time computation at the edge of chaos in recurrent neural networks. *Neural Computation*, 16(7):1413–1436, July 2004. ISSN 0899-7667. doi: 10.1162/089976604323057443.
- [2] Youngmin Cho and Lawrence K. Saul. Kernel methods for deep learning. In *Advances in neural information processing systems*, pages 342–350, 2009. URL <http://papers.nips.cc/paper/3628-kernel-methods-for-deep-learning>.
- [3] Amit Daniely, Roy Frostig, and Yoram Singer. Toward Deeper Understanding of Neural Networks: The Power of Initialization and a Dual View on Expressivity. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 2253–2261. Curran Associates, Inc., 2016. URL <http://papers.nips.cc/paper/6427-toward-deeper-understanding-of-neural-networks-the-power-of-initialization-and-a-dual-pdf>.
- [4] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *PMLR*, pages 249–256, March 2010. URL <http://proceedings.mlr.press/v9/glorot10a.html>.
- [5] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034, 2015. URL [http://www.cv-foundation.org/openaccess/content\\_iccv\\_2015/html/He\\_Delving\\_Deep\\_into\\_ICCV\\_2015\\_paper.html](http://www.cv-foundation.org/openaccess/content_iccv_2015/html/He_Delving_Deep_into_ICCV_2015_paper.html).
- [6] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image Recognition. pages 770–778, 2016. URL [https://www.cv-foundation.org/openaccess/content\\_cvpr\\_2016/html/He\\_Deep\\_Residual\\_Learning\\_CVPR\\_2016\\_paper.html](https://www.cv-foundation.org/openaccess/content_cvpr_2016/html/He_Deep_Residual_Learning_CVPR_2016_paper.html).
- [7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In *European Conference on Computer Vision*, pages 630–645. Springer, 2016.
- [8] Timothy P. Lillicrap, Daniel Cownden, Douglas B. Tweed, and Colin J. Akerman. Random synaptic feedback weights support error backpropagation for deep learning. *Nature Communications*, 7:ncomms13276, November 2016. ISSN 2041-1723. doi: 10.1038/ncomms13276. URL <https://www.nature.com/articles/ncomms13276>.
- [9] Ben Poole, Subhaneil Lahiri, Maithreyi Raghu, Jascha Sohl-Dickstein, and Surya Ganguli. Exponential expressivity in deep neural networks through transient chaos. In *Advances In Neural Information Processing Systems*, pages 3360–3368, 2016.
- [10] Maithra Raghu, Ben Poole, Jon Kleinberg, Surya Ganguli, and Jascha Sohl-Dickstein. On the expressive power of deep neural networks. *arXiv:1606.05336 [cs, stat]*, June 2016. URL <http://arxiv.org/abs/1606.05336>. arXiv: 1606.05336.
- [11] Samuel S. Schoenholz, Justin Gilmer, Surya Ganguli, and Jascha Sohl-Dickstein. Deep Information Propagation. 2017. URL <https://openreview.net/pdf?id=H1W1UN9gg>.

## Notes

<sup>1</sup>Under simplified conditions, Daniely et al. [3] showed that there exists a fixed point for any “well-behaved” activation function in a feedforward net. However, this result does not apply to architectures with residual connections.

<sup>2</sup>Note that in practice, to avoid the diverging gradient  $\dot{\psi}_\alpha(x) \rightarrow \infty$  as  $x \rightarrow 0$ , we can use a tempered version  $\Psi_\alpha(x)$  of  $\alpha$ -ReLU, defined by  $\Psi_\alpha(x) = (x + \epsilon)^\alpha - \epsilon^\alpha$  on  $x > 0$  and 0 otherwise, for some small  $\epsilon > 0$ . The conclusions of this paper on  $\psi_\alpha$  should hold similarly for  $\Psi_\alpha$  as well.

<sup>3</sup>Daniely et al. [3] called the version of  $W\phi$  with fixed  $\rho = 1$  the “dual function” of  $\phi$ .

<sup>4</sup>A more natural visualization is to graph  $\mathbf{e}^{(l)} - \mathbf{e}^*$  versus  $l^{-\delta^*}$ , but because of floating point precision,  $\mathbf{e}^{(l)} - \mathbf{e}^*$  doesn’t converge to 0, but a small number close to 0, so that the log-log plot wouldn’t look like what is expected.

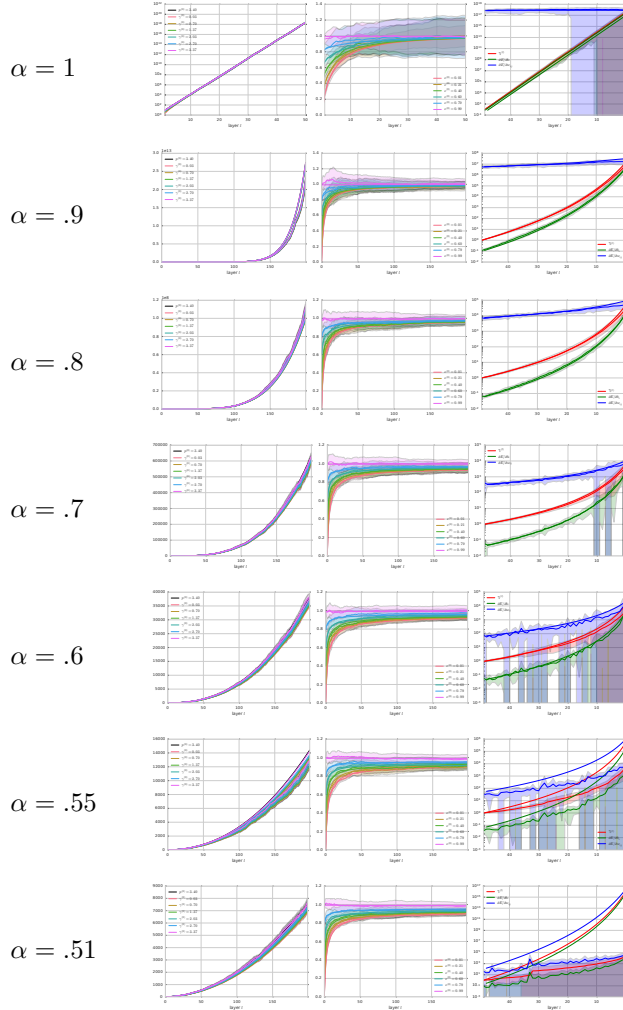
<sup>5</sup>Our derivations actually apply to all  $\alpha \in (\frac{1}{2}, 1]$ , where at  $\alpha = \frac{1}{2}$ , the expected norm of the gradient diverges within our mean field formalism. However, at  $\alpha \leq \frac{3}{4}$ , the variance of the gradient already diverges (Thm B.19), so we cannot expect the empirical values to agree with our theoretical predictions. But in fact, empirically our theoretical predictions seem to form an upper bound on the gradient norms (see Fig. A.1).

<sup>6</sup>the contour for  $\mathbf{p}^{(l)}$  is similar, but its slopes are slightly off from the heatmap contours.

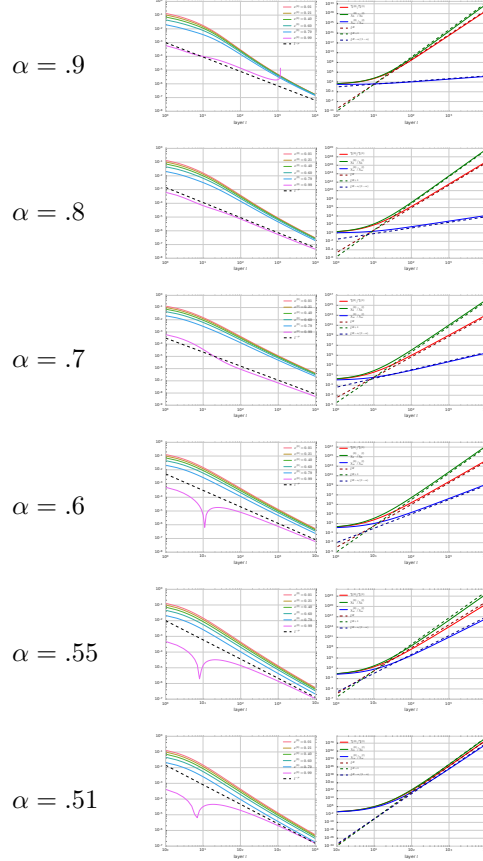
# Appendices

## A Additional Figures

In figures appearing in the appendix,  $\Upsilon$  means  $\chi$  (due to legacy reasons).



**Figure A.1:** Empirical vs theoretical dynamics for  $\mathbf{p}^{(l)}$ ,  $\mathbf{e}^{(l)}$ , and different gradient quantities for  $\alpha$ -ReLU, with format similar to Fig. 1. We refer to each figure on each row from left to right as (a), (b), and (c). Note that in the  $\alpha = 1$  case, figure (a) ( $\mathbf{p}^{(l)}$  and  $\gamma^{(l)}$  for different initial values) has log scale y-axis and (a) and (b) have x-axis ranging from 1 to 50, while for other  $\alpha$ , (a) has normal y-axis and (a) and (b) have x-axis ranging from 1 to 200. We do so because the norm of the activation vector in a typical ReLU resnet blows up into NaN at around layer 90, while this is not a problem for  $\alpha < 1$ . Our theoretical predictions track the average of empirical values closely for forward quantities  $\mathbf{p}^{(l)}$ ,  $\gamma^{(l)}$ , and  $\mathbf{e}^{(l)}$  for all  $\alpha$ , but variance is extremely large for  $\mathbf{e}^{(l)}$  at  $\alpha = 1$ ; it also predicts the average gradient norm accurately for  $\alpha = 1$  to  $\alpha = .7$  (despite the fact that we should not expect so for  $\alpha \leq .75$  due to exploding variance (Thm B.19)), although variance is large for  $\alpha = 1$  at earlier layers (i.e. later layers w.r.t backpropagation). However it *consistently and significantly overestimates* the average gradient norm for  $\alpha = .6$  to  $\alpha = .5$ , where the variance is so large that one standard deviation below the mean results in negative values. All plots are made with parameters  $\sigma_v^2 = 1.5$ ,  $\sigma_a^2 = .5$ ,  $\sigma_w^2 = 1.69$ ,  $\sigma_b^2 = .49$ ; only  $\alpha$  is varied. All figures exhibit smooth curves, which are theoretical estimates, and irregular curves with shades around them, which indicate empirical means and standard deviations (both of which taken in regular scale, not log scale). For each  $\alpha$ , figures (a) and (b) are made with 20 runs of resnets of width 1000. (c) is made with 25 runs of resnets of width 250.



**Figure A.2:** We verify the exponents of the forward and backward dynamics for  $\alpha$ -ReLU FRN. For each row, the figures are labeled (a) and (b) from left to right. The format is the same as in Fig. C.17. All figures are in log-log scale. (a) We exhibit our theoretical dynamics of the cosine distance  $\mathbf{e}^{(l)}$  based on the recurrences Thm B.8 and Thm B.10 for different initial conditions  $\mathbf{e}^{(0)}$ . We draw  $|\mathbf{e}^{(l)} - \mathbf{e}^{(l-1)}|$  for each of these dynamics in colored solid lines. We predict that each dynamic is  $\tilde{\Theta}(l^{-\mu})$ , where  $\mu = (1 - \mathbb{J}_\alpha(\mathbf{e}^*)) / (1 - \alpha)$ , and the dashed line gives  $l^{-\mu-1}$  (Thm B.18), shifted vertically to better compare the slope in log scale (i.e. the exponent of the polynomial dynamics). (See footnote 4 for why we plot the dynamics this way). We see that the our asymptotic prediction is very accurate for the sequence of  $\mathbf{e}^{(l)}$  that starts with  $\mathbf{e}^{(0)} = 0.99$ , the closest to  $\mathbf{e}^*$  for each  $\alpha$ , while other lines only slowly converge to the same exponent (which is the slope in the log-log plot). This is to be expected based on the proof of Thm B.18. For  $\alpha = .9$ , the  $\mathbf{e}^{(0)} = .99$  line upticks at around  $10^3$  and then turn into NaNs due to numerical instability. (b) Colored lines are  $\bullet^{(0)}/\bullet^{(l)}$  for  $\bullet = \chi, \chi_b, \chi_w$  (we are not taking logs in addition to plotting in log-log scale like in Fig. C.15). The dashed lines are our asymptotic predictions for the dynamics with corresponding colors, based on Thm B.21, again shifted appropriately to easily compare slope visually. We see that for every alpha our asymptotic predictions are highly accurate. For both (a) and (b), we did not show  $\alpha = 1$  case as ReLU FRN runs into numerical issues quickly (i.e. with even for 100 layers) because of exponential explosions in  $\mathbf{p}^{(l)}$  and  $\chi^{(l)}$  as predicted by Thms B.16 and B.20, so we cannot expect to empirically verify the precise predicted asymptotics. All plots are made with parameters  $\sigma_v^2 = 1.5, \sigma_a^2 = .5, \sigma_w^2 = 1.69, \sigma_b^2 = .49$ ; only  $\alpha$  is varied.

**Table A.1:** Glossary of Symbols. “Mean normalized” is abbreviated “m.n.”

Symbol	Meaning	Ref
$\sigma_\bullet$	standard deviation of trainable parameter $\bullet$	
$x^{(l)}$	activation vector/input vector	
$h^{(l)}$	hidden vector	
$N$	width (same across all layers)	
$\mathbf{p}^{(l)}$	m.n. squared length of activation vector $x^{(l)}$	3.3
$\mathbf{q}^{(l)}$	m.n. squared length of hidden vector $h^{(l)}$	3.3
$\gamma^{(l)}$	m.n. dot product $x^{(l)} \cdot x^{(l) \prime}$	3.4
$\lambda^{(l)}$	m.n. dot product $h^{(l)} \cdot h^{(l) \prime}$	3.4
$\mathbf{s}^{(l)}$	m.n. squared distance $\ x^{(l)} - x^{(l) \prime}\ ^2$	3.4
$\mathbf{e}^{(l)}$	cosine distance $\gamma^{(l)} / \sqrt{\mathbf{p}^{(l)} \mathbf{p}^{(l) \prime}}$	3.4
$\mathbf{e}^*$	limit value of $\mathbf{e}^{(l)}$ as $l \rightarrow \infty$	
$\mathbf{c}^{(l)}$	cosine distance $\lambda^{(l)} / \sqrt{\mathbf{q}^{(l)} \mathbf{q}^{(l) \prime}}$	3.4
$\chi^{(l)}$	m.n. gradient squared norm w.r.t. $x^{(l)}$	3.5
$\chi_\bullet^{(l)}$	m.n. gradient squared norm w.r.t. trainable parameter $\bullet$	3.5
$\phi$	variable nonlinearity $\mathbb{R} \rightarrow \mathbb{R}$	
$\psi_\alpha$	$\alpha$ -ReLU	5.2
$V$	variance integral transform	5.3
$W$	covariance integral transform	5.3
$\delta^*$	$\mathbf{e}^{(l)}$ converges like $\Theta(l^{-\delta^*})$ in tanh FRN	B.11
$\mathcal{A}$	leading coeff of $\log \chi^{(0)} / \chi^{(L)}$ in tanh FRN	B.13
$R$	$\log \chi^{(0)} / \chi^{(L)} \sim R \log L$ for $(\alpha < 1)$ -ReLU	B.20
$\mathbb{J}_\alpha$	kernel function of $\alpha$ -ReLU	C.30

## B A Listing of Main Theorems

### B.1 Tanh

#### B.1.1 Reduced Residual Network

**Lemma B.1.** Suppose  $\phi$  is antisymmetric. Then in an RRN,  $\mathbf{p}$  and  $\mathbf{q}$  satisfy the recurrence

$$\begin{aligned} \mathbf{q} &= \sigma_w^2 \underline{\mathbf{p}} + \sigma_b^2 \\ \mathbf{p} &= V\phi(\mathbf{q}) + \underline{\mathbf{p}}. \end{aligned}$$

**Theorem B.2.** Suppose  $\phi$  is tanh-like. Assume RRN architecture.

- If  $\sigma_w = 0$ , then  $\mathbf{p}^{(l)} = lV\phi(\sigma_b^2) + \mathbf{p}^{(0)}$  and  $\mathbf{q}^{(l)} = \sigma_b^2$ .
- If  $\sigma_w > 0$ ,  $\lim_{l \rightarrow \infty} \mathbf{p}^{(l)}/l = 1$  and  $\lim_{l \rightarrow \infty} \mathbf{q}^{(l)}/(\sigma_w^2 l) = 1$ . If  $\phi = \tanh$ , then we can obtain more terms of the asymptotic expansions:

$$\begin{aligned} \mathbf{p}^{(l)} &= l - 2C\sigma_w^{-1}l^{1/2} - C^2\sigma_w^{-2} \log l + O(1) \\ \mathbf{q}^{(l)} &= \sigma_w^2 l - 2C\sigma_w l^{1/2} - C^2 \log l + O(1) \end{aligned}$$

as  $l \rightarrow \infty$ , where  $C = \sqrt{2/\pi}$ .

**Theorem B.3.** Suppose  $\phi$  is antisymmetric. Then in an RRN,  $\lambda$  and  $\gamma$  satisfy the recurrence

$$\begin{aligned} \lambda &= \sigma_w^2 \underline{\gamma} + \sigma_b^2 \\ \gamma &= W\phi(\mathbf{q}, \lambda) + \underline{\gamma}. \end{aligned}$$

**Theorem B.4.** Suppose  $\phi$  is a tanh-like nonlinearity in an RRN. Assume  $\mathbf{e}^{(0)} < 1$ .

- If  $\sigma_w = 0$ , then  $\gamma^{(l)} = lW\phi(\sigma_b^2, \sigma_b^2) + \gamma^{(0)} = lV\phi(\sigma_b^2) + \gamma^{(0)}$  and  $\lambda^{(l)} = \sigma_b^2$ , so that  $\mathbf{e}^{(l)} \rightarrow 1$  and  $1 - \mathbf{e}^{(l)} = \Theta(l^{-1})$ . As a result,  $\mathbf{s}^{(l)} = \mathbf{p}^{(l)}(1 - \mathbf{e}^{(l)}) = \Theta(1)$ .

- If  $\sigma_w > 0$ , then  $\gamma^{(l)} = \check{\Theta}(l^{\frac{2}{\pi}})$ , and  $\mathbf{e}^{(l)} \rightarrow 0$  like  $\check{\Theta}(l^{\frac{2}{\pi}-1})$ . Thus  $\mathbf{s}^{(l)} = \Theta(\mathbf{p}^{(l)}) = \Theta(l)$ .

**Theorem B.5.** For any nonlinearity  $\phi$  in an RRN, under assumptions [Axiom 3.1](#) and [Axiom 3.2](#), whenever  $\dot{\phi}^2(\zeta)$  has finite variance for Gaussian variable  $\zeta$ ,

$$\underline{\chi} = (\sigma_w^2 \mathbf{V}\dot{\phi}(\mathbf{q}) + 1)\underline{\chi}, \quad \underline{\chi}_b = \underline{\chi} \mathbf{V}\dot{\phi}(\mathbf{q}), \quad \underline{\chi}_w = \underline{\chi} \mathbf{V}\dot{\phi}(\mathbf{q})\underline{\mathbf{p}}.$$

**Theorem B.6.** For  $\phi = \tanh$  in an RRN,

- If  $\sigma_w = 0$ ,  $\underline{\chi}^{(m)} = \underline{\chi}^{(l)}$  for all  $l, m$ .
- If  $\sigma_w > 0$ ,

$$\log(\underline{\chi}^{(m)}/\underline{\chi}^{(l)}) = \mathcal{A}(\sqrt{l} - \sqrt{m}) + \mathcal{B}(\log l - \log m) + O(1)$$

$$\text{where } \mathcal{A} = \frac{4}{3}\sqrt{\frac{2}{\pi}}\sigma_w \text{ and } \mathcal{B} = \frac{4}{3\pi} - \sigma_w^2 \frac{4}{9\pi}.$$

**Theorem B.7.** Suppose  $\phi = \tanh$ . Then in an RRN

- If  $\sigma_w = 0$ ,  $\underline{\chi}_b^{(l)} = \underline{\chi}^{(L)} \mathbf{V}\dot{\phi}(\sigma_b^2)$  and  $\underline{\chi}_w^{(l)} = \underline{\chi}^{(L)} \mathbf{V}\dot{\phi}(\sigma_b^2)((l-1)\mathbf{V}\phi(\sigma_b^2) + \mathbf{p}^{(0)})$ , where  $L$  is the last layer.
- If  $\sigma_w > 0$ ,

$$\log(\underline{\chi}_b^{(m)}/\underline{\chi}_b^{(l)}) = \mathcal{A}(\sqrt{l} - \sqrt{m}) + \mathcal{B}_b(\log l - \log m) + O(1)$$

$$\log(\underline{\chi}_w^{(m)}/\underline{\chi}_w^{(l)}) = \mathcal{A}(\sqrt{l} - \sqrt{m}) + \mathcal{B}_w(\log l - \log m) + O(1)$$

$$\text{where } \mathcal{A} = \frac{4}{3}\sqrt{\frac{2}{\pi}}\sigma_w \text{ (same as } \mathcal{A} \text{ in } \text{Thm B.6}) \text{ and } \mathcal{B}_b = \mathcal{B} + \frac{1}{2}, \mathcal{B}_w = \mathcal{B} - \frac{1}{2}, \text{ with } \mathcal{B} = \frac{4}{3\pi} - \sigma_w^2 \frac{4}{9\pi} \text{ (same as } \mathcal{B} \text{ in } \text{Thm B.6}).$$

### B.1.2 Full Residual Network

**Theorem B.8.** For any nonlinearity  $\phi$  in an FRN,

$$\mathbf{q} = \sigma_w^2 \underline{\mathbf{p}} + \sigma_b^2$$

$$\underline{\mathbf{p}} = \sigma_v^2 \mathbf{V}\phi(\mathbf{q}) + \sigma_a^2 + \underline{\mathbf{p}}$$

**Theorem B.9.** Suppose  $\phi$  is tanh-like. Assume the FRN architecture.

- If  $\sigma_w = 0$ , then  $\mathbf{p}^{(l)} = (\sigma_v^2 \mathbf{V}\phi(\sigma_b^2) + \sigma_a^2)l + \mathbf{p}^{(0)}$ , and  $\mathbf{q}^{(l)} = \sigma_b^2$ .
- If  $\sigma_w > 0$ , then  $\mathbf{p}^{(l)} = b_0 l + b_1 l^{1/2} + b_2 \log l + O(1)$ , where

$$b_0 = \sigma_v^2 + \sigma_a^2$$

$$b_1 = \frac{-2C\sigma_v^2\sigma_w^{-1}}{\sqrt{\sigma_v^2 + \sigma_a^2}}$$

$$b_2 = \frac{-C^2\sigma_v^4\sigma_w^{-2}}{(\sigma_v^2 + \sigma_a^2)^2}$$

$$\text{and } C = \sqrt{\frac{2}{\pi}}. \text{ Additionally, } \mathbf{q}^{(l)} = \sigma_w^2 b_0 l + \sigma_w^2 b_1 l^{1/2} + \sigma_w^2 b_2 \log l + O(1).$$

**Theorem B.10.** For any nonlinearity  $\phi$ , in an FRN

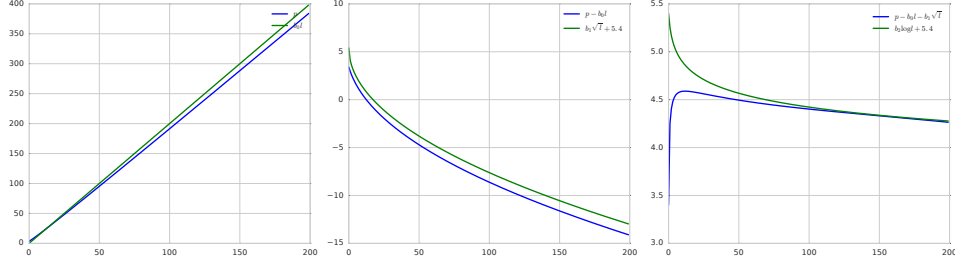
$$\underline{\lambda} = \sigma_w^2 \underline{\gamma} + \sigma_b^2$$

$$\underline{\gamma} = \sigma_v^2 \mathbf{W}\phi(\mathbf{q}, \underline{\lambda}) + \sigma_a^2 + \underline{\gamma}$$

**Theorem B.11.** Assume  $\phi = \tanh$  in an FRN. Suppose  $\mathbf{e}^{(0)} < 1$ .

- If  $\sigma_w = 0$ , then  $\underline{\lambda}^{(l)} = \sigma_b^2$  and  $\underline{\gamma}^{(l)} = l(\sigma_v^2 \mathbf{W}\phi(\sigma_b^2, \sigma_b^2) + \sigma_a^2) + \underline{\gamma}^{(0)} = l(\sigma_v^2 \mathbf{V}\phi(\sigma_b^2) + \sigma_a^2) + \underline{\gamma}^{(0)}$ . Thus  $\mathbf{e}^{(l)} \rightarrow 1$  and  $1 - \mathbf{e}^{(l)} = \Theta(l^{-1})$ . As a result,  $\mathbf{s}^{(l)} = \mathbf{p}^{(l)}(1 - \mathbf{e}^{(l)}) = \Theta(1)$ .





**Figure B.3:** Empirical verification of **Thm B.9**.

- If  $\sigma_w > 0$ , then  $\mathbf{e}^{(l)}$  converges to the unique fixed point  $\mathbf{e}^* \neq 1$  determined by the equation

$$\mathbf{e}^* = \frac{1}{\sigma_v^2 + \sigma_a^2} \left[ \sigma_v^2 \frac{2}{\pi} \arcsin(\mathbf{e}^*) + \sigma_a^2 \right].$$

Furthermore,  $\mathbf{e}^{(l)}$  converges to  $\mathbf{e}^*$  polynomially:  $|\mathbf{e}^{(l)} - \mathbf{e}^*|$  is  $\tilde{\Theta}(l^{-\delta^*})$ , where

$$\delta^* := 1 - \frac{2}{\pi} \frac{1}{\sqrt{1 - (\mathbf{e}^*)^2}} \frac{\sigma_v^2}{\sigma_v^2 + \sigma_a^2} \in \left[ \frac{2}{\pi} - 1, \frac{1}{2} \right)$$

Since  $\mathbf{e}^* < 1$ ,  $\mathbf{s}^{(l)} = \Theta(\mathbf{p}^{(l)}) = \Theta(l)$ .

**Theorem B.12.** For any nonlinearity  $\phi$  in an FRN, under assumptions **Axiom 3.1** and **Axiom 3.2**, whenever  $\dot{\phi}(\zeta)^2$  has finite variance for Gaussian variable  $\zeta$ ,

$$\begin{aligned} \underline{\chi} &= (\sigma_v^2 \sigma_w^2 \mathbf{V} \dot{\phi}(\mathbf{q}) + 1) \underline{\chi}, & \chi_b &= \sigma_v^2 \chi \mathbf{V} \dot{\phi}(\mathbf{q}), \\ \chi_w &= \sigma_v^2 \chi \mathbf{V} \dot{\phi}(\mathbf{q}) \underline{\mathbf{p}}, & \chi_v &= \chi \mathbf{V} \phi(\mathbf{q}), & \chi_a &= \chi \end{aligned}$$

**Theorem B.13.** Assume  $\phi = \tanh$  in an FRN.

- If  $\sigma_w = 0$ ,  $\chi^{(m)} = \chi^{(l)}$  for all  $l, m$ .
- If  $\sigma_w > 0$ , then for  $l \geq m \geq 0$ ,

$$\log(\chi^{(m)}/\chi^{(l)}) = \mathcal{A}(\sqrt{l} - \sqrt{m}) + \mathcal{B}(\log l - \log m) + O(1)$$

where

$$\begin{aligned} \mathcal{A} &= \frac{4}{3} \sqrt{\frac{2}{\pi}} \frac{\sigma_v^2 \sigma_w}{\sqrt{\sigma_v^2 + \sigma_a^2}} \\ \mathcal{B} &= \frac{4}{9\pi} \frac{\sigma_v^4}{\sigma_v^2 + \sigma_a^2} \left( \frac{3}{\sigma_v^2 + \sigma_a^2} - \sigma_w^2 \right) \end{aligned}$$

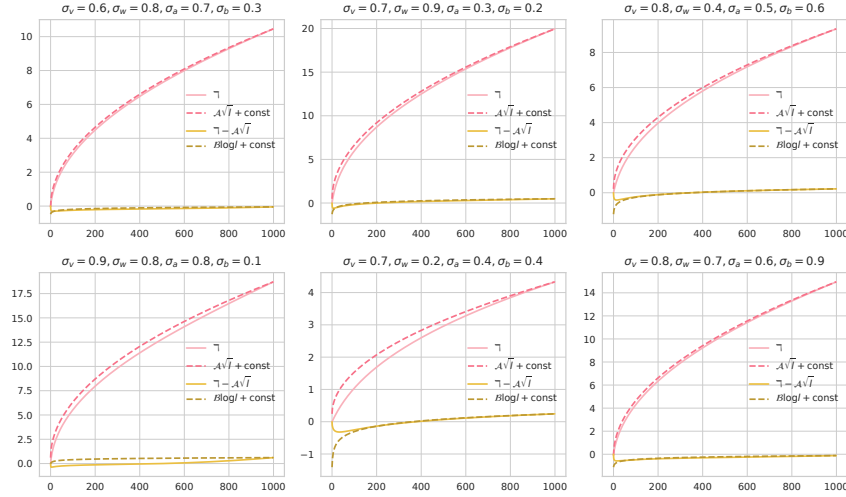
**Fig. B.4** shows empirical verification of the asymptotic expansion of  $\chi$  for various values of  $\sigma_\bullet$ s.

**Theorem B.14.** Suppose  $\phi = \tanh$  in an FRN.

- If  $\sigma_w = 0$ , then

$$\begin{aligned} \chi_b^{(l)} &= \sigma_v^2 \chi^{(L)} \mathbf{V} \dot{\phi}(\sigma_b^2) \\ \chi_w^{(l)} &= \sigma_v^2 \chi^{(L)} \mathbf{V} \dot{\phi}(\sigma_b^2) ((\sigma_v^2 \mathbf{V} \phi(\sigma_b^2) + \sigma_a^2)(l-1) + \mathbf{p}^{(0)}) \\ \chi_v^{(l)} &= \chi^{(L)} \mathbf{V} \phi(\sigma_b^2) \\ \chi_a^{(l)} &= \chi^{(L)}. \end{aligned}$$

- If  $\sigma_w > 0$ , then for  $l \geq m \geq 0$ ,



**Figure B.4:** Empirical verification of the asymptotic expansion of  $\chi$  for various values of  $\sigma_{\bullet}$ s. Note that we have chosen all small values for  $\sigma_{\bullet}$ s. For larger values, the constant term in [Thm B.13](#) begins to dominate (primarily because of the expansion  $\log(1+x) = x + \Theta(x^2)$  (has large  $\Theta$  term when  $x$  is large), and  $\chi$  behaves more like  $\exp(\Theta(l))$  up to depth 1000.

$$\begin{aligned} \log(\chi_b^{(m)}/\chi_b^{(l)}) &= \mathcal{A}(\sqrt{l} - \sqrt{m}) + \mathcal{B}_b(\log l - \log m) + O(1) \\ \log(\chi_w^{(m)}/\chi_w^{(l)}) &= \mathcal{A}(\sqrt{l} - \sqrt{m}) + \mathcal{B}_w(\log l - \log m) + O(1) \\ \log(\chi_a^{(m)}/\chi_a^{(l)}) &= \mathcal{A}(\sqrt{l} - \sqrt{m}) + \mathcal{B}(\log l - \log m) + O(1) \\ \log(\chi_v^{(m)}/\chi_v^{(l)}) &= \mathcal{A}(\sqrt{l} - \sqrt{m}) + \mathcal{B}(\log l - \log m) + O(1) \end{aligned}$$

where  $\mathcal{A} = \frac{4}{3} \sqrt{\frac{2}{\pi}} \frac{\sigma_v^2 \sigma_w}{\sqrt{\sigma_v^2 + \sigma_a^2}}$  and  $\mathcal{B} = \frac{4}{9\pi} \frac{\sigma_v^4}{\sigma_v^2 + \sigma_a^2} \left( \frac{3}{\sigma_v^2 + \sigma_a^2} - \sigma_w^2 \right)$  are as in [Thm B.13](#) and  $\mathcal{B}_b = \mathcal{B} + \frac{1}{2}$  and  $\mathcal{B}_w = \mathcal{B} - \frac{1}{2}$ .

## B.2 $\alpha$ -ReLU

**Lemma B.15.** If  $\alpha > -\frac{1}{2}$ , then  $\mathbb{V}\psi_\alpha(q) = c_\alpha q^\alpha$ , where  $c_\alpha = \frac{1}{\sqrt{\pi}} 2^{\alpha-1} \Gamma(\alpha + \frac{1}{2})$ .

Note that if  $\alpha \leq -\frac{1}{2}$ , then  $\mathbb{V}\psi_\alpha(q)$  is not defined (its defining integral does not converge).

### B.2.1 Full Residual Network

By [Thm B.8](#) and [Lemma B.15](#), we have the length recurrences

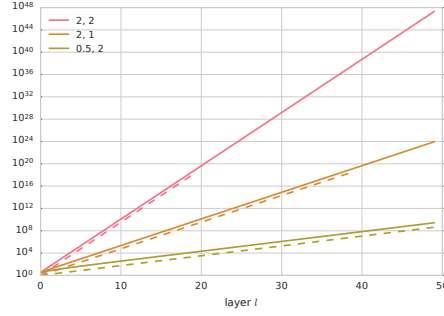
$$\begin{aligned} \mathbf{q} &= \sigma_w^2 \underline{\mathbf{p}} + \sigma_b^2 \\ \mathbf{p} &= \sigma_v^2 c_\alpha \mathbf{q}^\alpha + \sigma_a^2 + \underline{\mathbf{p}} \end{aligned}$$

**Theorem B.16.** Suppose we have the nonlinearity  $\phi = \psi_\alpha$ . The in an FRN: If  $\alpha = 1$ , then  $\mathbf{p}^{(l)} = \Theta((1 + \sigma_v^2 \sigma_w^2 / 2)^l)$ , with the hidden constant depending on the initial condition. If  $0 < \alpha < 1$ , then  $\mathbf{p}^{(l)} = \Theta(l^{\frac{1}{1-\alpha}})$ . More precisely,  $\lim_{l \rightarrow \infty} \mathbf{p} / l^{\frac{1}{1-\alpha}} = [\sigma_v^2 \sigma_w^2 c_\alpha (1 - \alpha)]^{\frac{1}{1-\alpha}}$ .

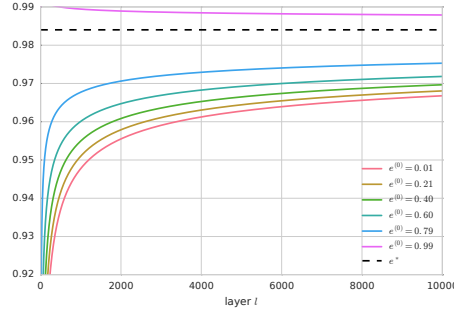
[Fig. B.5](#) empirically verifies the asymptotics for  $\alpha = 1$  for various  $\sigma_v$  and  $\sigma_w$ .

Similarly, by [Thm B.10](#), if  $\mathbf{q} = \mathbf{q}'$ , then

$$\begin{aligned} \lambda &= \sigma_w^2 \underline{\gamma} + \sigma_b^2 \\ \gamma &= \sigma_v^2 \mathbf{q}'^\alpha \mathbb{W}\psi_\alpha(1, \mathbf{c}) + \sigma_a^2 + \underline{\gamma} \end{aligned}$$



**Figure B.5:** Verification of the exponential asymptotics of  $\mathbf{p}^{(l)}$  when  $\alpha = 1$ . The lines of each color correspond to different  $(\sigma_w, \sigma_v)$  pairs, which are given in the legend. The solid lines are given by the recurrences [Thm B.8](#), and the dashed lines are given by our asymptotics  $(1 + \sigma_v^2 \sigma_w^2 / 2)^l$  ([Thm B.16](#)). Note that the y-axis is in log-scale.



**Figure B.6:** Verification of fixed point  $\mathbf{e}^*$  in [Thm B.18](#) for  $\alpha = .6$ . Different colors correspond to different initial conditions  $\mathbf{e}^{(0)}$ , and the dashed line gives the fixed point.

**Theorem B.17.** Suppose  $\phi = \psi_1$ . Then in an FRN,  $\mathbf{e}^{(l)} \rightarrow 1$  and  $1 - \mathbf{e}^{(l)} \sim [\frac{1}{4} \sigma_v^2 \sigma_w^2 B^{-1} U l]^{-2}$  for  $B = 1 + \sigma_v^2 \sigma_w^2 / 2$  and  $U = \frac{2\sqrt{2}}{3\pi}$ . As a result,  $\mathbf{s}^{(l)} = (1 - \mathbf{e}^{(l)}) \mathbf{p}^{(l)} = \Theta(l^{-2} \exp(\Theta(l))) = \exp(\Theta(l))$ .

**Theorem B.18.** Suppose  $\phi = \psi_\alpha$  for  $0 < \alpha < 1$  in an FRN. Then  $\mathbf{e}$  converges to the unique nonunit fixed point  $\mathbf{e}^*$  of  $\mathbb{J}_\alpha$ , and  $|\mathbf{e}^* - \mathbf{e}^{(l)}|$  is  $\Theta(l^{-\mu})$ , where  $\mu = (1 - \mathbb{J}_\alpha(\mathbf{e}^*)) / (1 - \alpha)$ . Additionally,  $\mathbf{s}^{(l)} = \Theta(\mathbf{p}^{(l)}) = \Theta(l^{1/(1-\alpha)})$ .

[Fig. B.6](#) verifies empirically that  $\mathbf{e}^*$  is indeed the fixed point of  $\mathbf{e}^{(l)}$ . [Fig. A.2](#) verifies empirically the convergence rate  $l^{-\mu}$ . [Fig. B.7](#) plots  $\mathbb{J}_\alpha(\mathbf{e}^*)$  and  $\mu$  versus  $\alpha$ . It certainly looks like  $\mu = \frac{1}{2}(1 - \alpha)$ , but we have no proof for it. Based on this conjecture, we see there is a “discontinuity” of  $\mu$  at  $\alpha = 1$ :  $\mu \rightarrow 0$  as  $\alpha \rightarrow 1$ , but for  $\alpha = 1$ , the actual convergence dynamics has exponent  $-2$  by [Thm B.17](#).

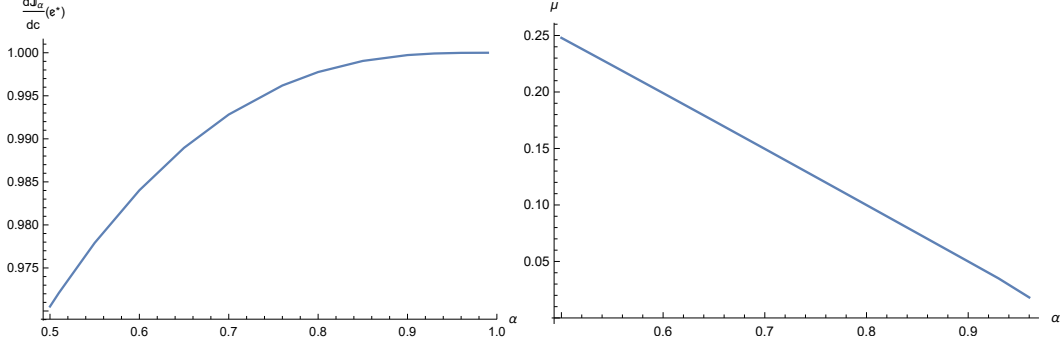
Because of the following theorem, we cannot expect the equations of [Thm B.12](#) to hold for  $\alpha \leq \frac{3}{4}$ .

**Theorem B.19.** Suppose we have the nonlinearity  $\psi_\alpha$  in an FRN.  $\text{Var}(\psi_\alpha(\zeta)^2)$  diverges for any Gaussian variable  $\zeta$  with mean 0 if  $\alpha \leq \frac{3}{4}$  but is finite if  $\alpha > \frac{3}{4}$ .

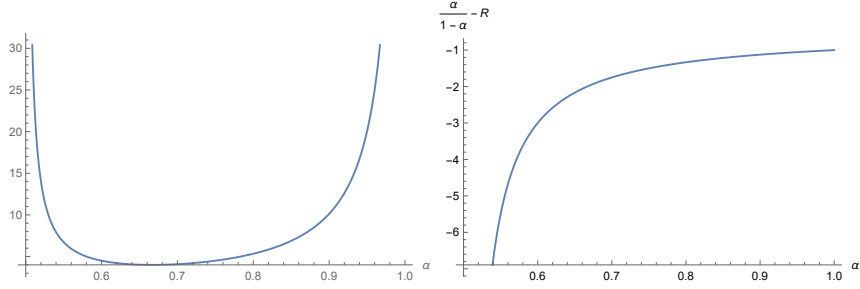
**Theorem B.20.** Suppose we have the nonlinearity  $\psi_\alpha$  in an FRN. If  $\alpha = 1$ , then  $\chi^{(l-m)} = \chi^{(l)} (\frac{1}{2} \sigma_v^2 \sigma_w^2 + 1)^m$ . If  $\alpha \in (\frac{3}{4}, 1)$ , then  $\chi^{(l-m)} = \Theta(1) \chi^{(l)} (l / (l - m))^R$  for  $R = \frac{\alpha^2}{(1-\alpha)(2\alpha-1)}$ , where the constants in  $\Theta(1)$  do not depend on  $l$  or  $m$ .

This exponent  $\frac{\alpha^2}{(1-\alpha)(2\alpha-1)}$  is minimized at  $\alpha = \frac{3}{4}$  on  $\alpha \in (3/4, 1)$ , where the value is  $\frac{9}{2}$  (and at  $\alpha = \frac{2}{3}$  on  $\alpha \in (1/2, 1)$ , where the value achieved is 4) ([Fig. B.8\(a\)](#)).

As a corollary,



**Figure B.7:** (a) A plot of  $\dot{J}_\alpha(\mathbf{e}^*)$  versus  $\alpha$ . (b) A plot of the exponent  $\mu$  of the dynamics of  $|\mathbf{e}^* - \mathbf{e}^{(l)}|$  (see Thm B.18)



**Figure B.8:** (a) The exponent of the polynomial gradient dynamics with respect to  $\alpha$ -ReLU versus  $\alpha$ . (b) The exponent of the dynamics of  $\chi_v$  and  $\chi_w$ .

**Theorem B.21.** If  $\phi = \psi_1$  in an FRN, then for  $l \geq m \geq 0$ ,

$$\begin{aligned} \chi_b^{(l-m)} &= \Theta(1)\chi^{(l)}B^m, & \chi_w^{(l-m)} &= \Theta(1)\chi^{(l)}B^l, \\ \chi_v^{(l-m)} &= \Theta(1)\chi^{(l)}B^l, & \chi_a^{(l-m)} &= \Theta(1)\chi^{(l)}B^m. \end{aligned}$$

where  $B = 1 + \sigma_v^2\sigma_w^2/2$ .

If  $\phi = \psi_\alpha$  in an FRN, for  $\alpha < 1$ , then for  $l \geq m \geq 0$ ,

$$\begin{aligned} \chi_b^{(l-m)} &= \Theta(1)\chi^{(l)}l^R(l-m)^{-R-1}, & \chi_w^{(l-m)} &= \Theta(1)\chi^{(l)}l^R(l-m)^{\frac{\alpha}{1-\alpha}-R}, \\ \chi_v^{(l-m)} &= \Theta(1)\chi^{(l)}l^R(l-m)^{\frac{\alpha}{1-\alpha}-R}, & \chi_a^{(l-m)} &= \Theta(1)\chi^{(l)}(l/(l-m))^R. \end{aligned}$$

Fig. A.2 verifies the backward asymptotic dynamics empirically for different  $\alpha < 1$ . Fig. B.8(b) graphs the exponent  $\frac{\alpha}{1-\alpha} - R$  in terms of  $\alpha$ . We see that on  $[0.5, 1]$ , the maximum of this exponent is at  $\alpha = 1$ .

## C Proofs

A brief note about notation: We use  $\sim$  to denote both how a random variable is sampled (ex:  $x \sim \mathcal{N}(0, 1)$  for a Gaussian  $x$ ) and how a function behaves asymptotically, i.e.  $f(x) \sim g(x)$  as  $x \rightarrow a$  iff  $\lim_{x \rightarrow a} f(x)/g(x) = 1$ . Context should be enough to differentiate between these two cases. We in addition use  $\simeq$  to denote asymptotic expansion. For example, if  $\{\alpha_i\}_{i \geq 0}$  is a sequence of strictly decreasing reals and  $\{\beta_i\}_{i \geq 0}$  is a sequence of nonzero reals, then

$$f(x) \simeq \sum_{i \geq 0} \beta_i (x - \xi)^{\alpha_i}$$

means that as  $x \rightarrow \xi$ ,  $f(x) - \sum_{i=0}^N \beta_i (x - \xi)^{\alpha_i} = \Theta((x - \xi)^{\alpha_{N+1}})$ .

## C.1 Preliminary Lemmas

**Lemma C.1.** *We have*

$$\frac{\sigma_w^2 \gamma + \sigma_b^2}{\sigma_w^2 \mathbf{p} + \sigma_b^2} = \mathbf{e}(1 + O(\gamma^{-1})).$$

regardless of whether  $\mathbf{e}^{(l)} = \gamma^{(l)}/\mathbf{p}^{(l)}$  converges.

But suppose  $\mathbf{e}^{(l)} = \gamma^{(l)}/\mathbf{p}^{(l)} \rightarrow \mathbf{e}^*$ . If  $\mathbf{e}^* < 1$ , then

$$\frac{\sigma_w^2 \gamma + \sigma_b^2}{\sigma_w^2 \mathbf{p} + \sigma_b^2} = \mathbf{e}(1 + \Theta(\gamma^{-1})).$$

If  $\mathbf{e}^* = 1$ , then

$$\frac{\sigma_w^2 \gamma + \sigma_b^2}{\sigma_w^2 \mathbf{p} + \sigma_b^2} = \mathbf{e}(1 + \Theta(\epsilon \mathbf{p}^{-1})),$$

where  $\epsilon = 1 - \mathbf{e}$ .

*Proof.* Write  $M = \sigma_b^2/\sigma_w^2$ .

$$\begin{aligned} \frac{\sigma_w^2 \gamma + \sigma_b^2}{\sigma_w^2 \mathbf{p} + \sigma_b^2} &= \mathbf{e} \left( 1 + \frac{1 + M\gamma^{-1}}{1 + M\mathbf{p}^{-1}} \right) \\ &= \mathbf{e} \left( 1 + M(\gamma^{-1} - \mathbf{p}^{-1}) + O(\mathbf{p}^{-1}(\gamma^{-1} - \mathbf{p}^{-1})) \right). \end{aligned}$$

In any situation,  $\gamma^{-1} - \mathbf{p}^{-1} = O(\gamma^{-1})$  because  $\gamma \leq \mathbf{p}$ , so this gives the first statement. If  $\mathbf{e}^*$  exists and  $\mathbf{e}^* < 1$ , then  $\gamma^{-1} - \mathbf{p}^{-1} = \Theta(\gamma^{-1})$ , which yields the second statement. If  $\mathbf{e}^*$  exists and  $\mathbf{e}^* = 1$ , then  $\gamma^{-1} - \mathbf{p}^{-1} = \mathbf{p}^{-1}((1 - \epsilon)^{-1} - 1) = \mathbf{p}^{-1}(\epsilon + O(\epsilon^2)) = \Theta(\epsilon \mathbf{p}^{-1})$ .  $\square$

For any function  $f$  that is  $(k+1)$ -times differentiable in a neighborhood of 0, we have the asymptotic expansion

$$f(z) = \sum_{n=0}^k \frac{d^n f}{dz^n}(0) \frac{z^n}{n!} + O(z^{k+1}), \text{ as } z \rightarrow 0.$$

Since

$$\frac{d^n}{d(1/q)^n} q^{1/2} \mathbb{V} \phi(q) \Big|_{q \rightarrow \infty} = \frac{(-1)^n}{2^n \sqrt{2\pi}} \int_{-\infty}^{\infty} \phi^2(z) z^{2n} dz$$

whenever the RHS is integrable, we have

**Lemma C.2.** *Suppose  $\phi^2(z)z^{2n}$  is integrable over  $z \in \mathbb{R}$  for all  $0 \leq n \leq N+1$ . Then  $\mathbb{V} \phi(q) = q^{-1/2}(\sum_{n=0}^N C_n q^{-n} + O(q^{-N-1}))$  as  $q \rightarrow \infty$ , where*

$$C_n := \frac{(-1)^n}{2^n n! \sqrt{2\pi}} \int_{-\infty}^{\infty} \phi^2(z) z^{2n} dz.$$

Note that  $\text{sech}^d(z) = \Theta(e^{-d|z|})$  for  $z \rightarrow \infty$  as long as  $d > 0$ , so that  $C_n$  from the above result converges when  $\phi = \text{sech}^d$ . Therefore

**Lemma C.3.** *Let  $d > 0$ . We have  $\mathbb{V} \text{sech}^d(q) \simeq q^{-1/2} \sum_{n \geq 0} C_n q^{-n}$ , where*

$$C_n := \frac{(-1)^n}{2^n n! \sqrt{2\pi}} \int_{-\infty}^{\infty} \text{sech}^{2d}(z) z^{2n} dz.$$

As corollaries, we obtain the following asymptotics.

**Lemma C.4.**  $\mathbb{V} \tanh(q) = \frac{2}{3} \sqrt{\frac{2}{\pi}} q^{-1/2} + \Theta(q^{-3/2})$  as  $q \rightarrow \infty$ .

*Proof.* Use **Lemma C.3** along with the fact that  $\dot{\tanh}(z) = \text{sech}^2(z)$  and  $\int \text{sech}^4 z dz = \frac{2}{3} \tanh z + \frac{1}{2} \text{sech}^2 z \tanh z$ .  $\square$

**Lemma C.5.**  $1 - V \tanh(q) = \sqrt{\frac{2}{\pi}} q^{-1/2} + \Theta(q^{-3/2})$  as  $q \rightarrow \infty$ .

*Proof.* Use **Lemma C.3** along with the fact that  $1 - \tanh^2(z) = \operatorname{sech}^2(z)$  and  $\int \operatorname{sech}^2 z \, dz = \tanh z$ .  $\square$

**Lemma C.6.**  $\operatorname{sech}^2(t) \geq \exp(-t^2)$  for all  $t$ , with equality iff  $t = 0$ .

*Proof.* The lower bound is equivalent to

$$2 \geq e^{t-t^2/2} + e^{-t-t^2/2}$$

The RHS has derivative  $(1-t)e^{t-t^2/2} - (1+t)e^{-t-t^2/2}$ . This is 0 iff

$$\frac{1-t}{1+t} = e^{-2t}$$

which has a solution 0 and in general can only have solution  $t \in (-1, 1)$  (by considering the sign of the LHS). Since each side is analytic in  $t \in (-1, 1)$ , we expand

$$\begin{aligned} \log \frac{1-t}{1+t} &= \log e^{-2t} \\ \log(1-t) - \log(1+t) &= -2t \\ (-t - t^2 - \dots) - (t - t^2 + \dots) &= -2t \\ -2t - 2t^3 - \dots &= -2t \end{aligned}$$

which shows that the only solution is  $t = 0$ . A simple plot shows that  $t = 0$  is a maximum, where the bound in question achieves equality.  $\square$

**Lemma C.7.** Suppose  $\phi = \tanh$ . Then  $V\dot{\phi}(q) \geq \frac{1}{\sqrt{4q+1}}$ .

As a sanity check, **Lemma C.4** shows that  $V\dot{\phi}(q) \sim C_0 q^{1/2}$  where  $C_0 \approx .5319$ , which is above the .5 in this lemma.

*Proof.* By **Lemma C.6**,

$$\begin{aligned} V\dot{\phi}(q) &= \int d\mu(z) \dot{\phi}^2(\sqrt{q}z) \\ &\geq \frac{1}{\sqrt{2\pi}} \int dz \exp(-z^2/2 - 2qz^2) \\ &= \frac{1}{\sqrt{2\pi}} \int dz \exp(-(4q+1)z^2/2) \\ &= \frac{1}{\sqrt{4q+1}}. \end{aligned}$$

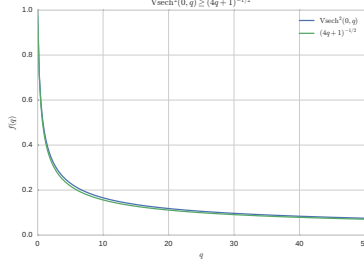
$\square$

**Fig. C.9** demonstrates **Lemma C.7**.

**Lemma C.8.** Let  $d \in \mathbb{R}$  and  $1 < M < N$  with  $N - M \in \mathbb{Z}^{\geq 0}$ . Set  $\Sigma(M, N, d) := \sum_{a=M}^N a^d$ . If we fix  $M$  and let  $N \rightarrow \infty$ ,

$$\Sigma(M, N, d) = \begin{cases} \Theta(1) & \text{if } d < -1 \\ \log N + O(1) & \text{if } d = -1 \\ \frac{N^{d+1}}{d+1} + O(1) & \text{if } -1 < d < 0 \\ N - M + 1 & \text{if } d = 0 \\ \frac{1}{d+1} N^{d+1} + \frac{1}{2} N^d + O(N^{\max(0, d-1)}) & \text{if } d > 0 \end{cases}$$





**Figure C.9:** Illustration of **Lemma C.7**:  $V\dot{\phi}(q)$  vs  $\frac{1}{\sqrt{4q+1}}$  for  $\phi = \tanh$ . This bound is very tight, and for most purposes,  $\frac{1}{\sqrt{4q+1}}$  can be taken as a good approximation of  $V\dot{\phi}(q)$ .

*Proof.* Consider the integrals  $A = \int_M^{N+1} a^d da$  and  $B = \int_{M-1}^N a^d da$ . They evaluate to  $A = \frac{1}{d+1}((N+1)^{d+1} - M^{d+1})$  and  $B = \frac{1}{d+1}(N^{d+1} - (M-1)^{d+1})$  when  $d \neq -1$  and to  $A = \log(N+1) - \log M$  and  $B = \log N - \log(M-1)$  when  $d = -1$ . When  $d \leq 0$ , we have  $A \leq B$  and  $\Sigma(M, N, d) \in [A, B]$ ; when  $d > 0$ ,  $B \leq A$  and  $\Sigma(M, N, d) \in [B, A]$ . Thus, as  $N \rightarrow \infty$  with  $M$  fixed, when  $d < -1$ ,  $\Sigma(M, N, d) = \Theta(1)$ ; when  $d = -1$ ,  $\Sigma(M, N, -1) = \log N + O(1)$ ; and when  $d > -1$ , we have  $\Sigma(M, N, d) = \frac{N^{d+1}}{d+1} + O(N^d)$ .

Now for  $a > 0$  and  $d > -1$  and  $d \neq 0, 1$ ,

$$\begin{aligned} \int_a^{a+1} z^d - a^d dz &= \frac{1}{d+1}((a+1)^{d+1} - a^{d+1}) \\ &= (a^d + \frac{d}{2}a^{d-1} + \dots) - a^d \\ &= \frac{d}{2}a^{d-1} + \Theta(a^{d-2}). \end{aligned}$$

where the hidden constants in  $\Theta$  depend only on  $d$  (and in fact this term vanishes if  $d = 1$ ). Thus

$$\begin{aligned} \Sigma(M, N, d) &= \int_M^{N+1} z^d dz - \sum_{a=M}^N [\frac{d}{2}a^{d-1} + \Theta(a^{d-2})] \\ &= \frac{1}{d+1}((N+1)^{d+1} - M^{d+1}) - \frac{d}{2}\Sigma(M, N, d-1) + \Theta(\Sigma(M, N, d-2)) \end{aligned}$$

If  $-1 < d < 0$ , then  $\Sigma(M, N, d-1) = \Theta(1)$ , so that  $\Sigma(M, N, d) = \frac{(N+1)^{d+1}}{d+1} + O(1) = \frac{N^{d+1}}{d+1} + O(1)$ . If  $d > 0$  and  $d \neq 1$ , then  $\Sigma(M, N, d-1) = \frac{N^d}{d}$ , so that

$$\begin{aligned} \Sigma(M, N, d) &= \frac{1}{d+1}N^{d+1} + N^d + \Theta(N^{\max(0, d-1)}) - \frac{1}{2}N^d + \Theta(\Sigma(M, N, d-2)) \\ &= \frac{1}{d+1}N^{d+1} + \frac{1}{2}N^d + O(N^{\max(0, d-1)}). \end{aligned}$$

□

We can obtain more terms in the expansion for higher  $d$  via the Euler-Maclaurin formula, but this suffices for our purposes.

## C.2 Dynamics Zoo

This section deduces the asymptotic behaviors of some sequences governed by recurrence equations. For the most part, the leading term of their asymptotic expansions is as one would expect from the corresponding differential equation. However, in some cases we need subleading terms for later results. They require slightly more nuanced reasoning. First we present a technical lemma.

**Lemma C.9.** *Let  $F : \mathbb{R} \times \mathbb{N} \rightarrow \mathbb{R}$  be a function such that for a subset  $U \subseteq \mathbb{R}$ , and for all  $z, z' \in U, z \geq z' \implies F(z, n) \geq F(z', n)$  for every  $n$ . Suppose sequences  $a^{(l)}, b^{(l)}, c^{(l)}$  satisfy*

- $a^{(l+1)} = F(a^{(l)}, l)$  for all  $l$ ;
- $b^{(l+1)} \leq F(b^{(l)}, l)$  for all  $l$  above a constant  $K_b$ .
- $c^{(l+1)} \geq F(c^{(l)}, l)$  for all  $l$  above a constant  $K_c$ .

and furthermore,  $a^{(l)}, b^{(l)}, c^{(l)}$  all fall into  $U$  for  $l$  above a constant  $K_U$ .

If for some  $m \geq \max(K_b, K_U)$ ,  $b^{(m)} \leq a^{(m)}$ , then  $b^{(l)} \leq a^{(l)}, \forall l \geq m$ . Similarly, if for some  $n \geq \max(K_c, K_U)$ ,  $c^{(n)} \geq a^{(n)}$ , then  $c^{(l)} \geq a^{(l)}, \forall l \geq n$ .

*Proof.* For the first claim:  $b^{(m)} \leq a^{(m)} \implies b^{(m+1)} \leq F(b^{(m)}, m) \leq F(a^{(m)}, m) = a^{(m+1)}$ . Here the last inequality used the monotonicity of  $F$ . Induction gives the desired result.

It's similar for the second claim, where the inductive step is  $c^{(m)} \geq a^{(m)} \implies c^{(m+1)} \geq F(c^{(m)}, m) \geq F(a^{(m)}, m) = a^{(m+1)}$ .  $\square$

**Lemma C.10.** Suppose  $\epsilon^{(l)}$  satisfies the recurrence

$$\epsilon^{(l)} = \epsilon^{(l-1)} \left(1 + \frac{\delta}{l^\beta}\right).$$

for some nonzero constant  $\delta \in \mathbb{R}$  independent of  $l$ .

- If  $\beta > 1$ , then  $\epsilon^{(l)} = \Theta(1)$ .
- If  $\beta = 1$ , then  $\epsilon^{(l)} = \Theta(l^\delta)$ .
- If  $0 < \beta < 1$ , then  $\epsilon^{(l)} = \exp\left(\frac{\delta}{1-\beta} l^{1-\beta} + \tilde{\Theta}(l^{\psi_1(1-2\beta)})\right)$ , where  $\psi_1(x) = \max(0, x)$  is the ReLU function.

*Proof.* We have

$$\begin{aligned} \log \epsilon^{(l)} &= \log \epsilon^{(l-1)} + \log(1 + \delta/l^\beta) \\ &= \log \epsilon^{(l-1)} + \delta/l^\beta + \Theta(\delta^2/l^{2\beta}) \end{aligned}$$

for large  $l$ . If  $\beta > 1$ , then  $\sum_l l^{-\beta}$  converges, and

$$\begin{aligned} \log \epsilon^{(l)} &= \log \epsilon^{(0)} - \Theta(1) \\ \epsilon^{(l)} &= \Theta(1). \end{aligned}$$

If  $\beta = 1$ , then

$$\begin{aligned} \log \epsilon^{(l)} &= \log \epsilon^{(0)} + \delta \log l + \Theta(1) \\ \epsilon^{(l)} &= \Theta(l^\delta). \end{aligned}$$

If  $\beta < 1$ , then

$$\begin{aligned} \log \epsilon^{(l)} &= \log \epsilon^{(0)} + \frac{\delta}{1-\beta} l^{1-\beta} + \tilde{\Theta}(l^{1-2\beta}) \\ \epsilon^{(l)} &= \exp\left(\frac{\delta}{1-\beta} l^{1-\beta} + \tilde{\Theta}(l^{\psi_1(1-2\beta)})\right). \end{aligned}$$

$\square$

**Lemma C.11.** Suppose  $\epsilon^{(l)} = Cl^{-\alpha} + \epsilon^{(l-1)}(1 + \delta/l^\beta)$  for  $\alpha \in \mathbb{R}$ ,  $C \neq 0$ , and  $\delta \neq 0$ . Then

- If  $\beta > 1$ , then
  - $\epsilon^{(l)} = \Theta(l^{1-\alpha})$  if  $\alpha \in (0, 1)$ ;
  - $\epsilon^{(l)} = \Theta(\log l)$  if  $\alpha = 1$ ;
  - $\epsilon^{(l)} = \Theta(1)$  if  $\alpha > 1$ .

- If  $\beta = 1$ , then
  - $\epsilon^{(l)} = \Theta(l^{\max(\delta, 1-\alpha)})$  if  $1 - \delta \neq \alpha$ .
  - $\epsilon^{(l)} = \Theta(l^\delta \log l)$  if  $1 - \delta = \alpha$ .

Furthermore, for  $\beta = -\delta = 1$ ,  $\epsilon^{(l)} \sim l^{-1}$  if  $\alpha > 2$ ,  $\epsilon^{(l)} \sim l^{1-\alpha}$  if  $\alpha < 2$ , and  $\epsilon^{(l)} \sim l^\delta \log l$  if  $\alpha = 2$ .

*Proof.* We can unwind the recurrence to get

$$\epsilon^{(l)} = \sum_{m=1}^l m^{-\alpha} \prod_{n=m+1}^l \left(1 + \frac{\delta}{n^\beta}\right) + \epsilon^{(0)} \prod_{n=1}^l \left(1 + \frac{\delta}{n^\beta}\right)$$

Suppose  $\beta > 1$ . By [Lemma C.10](#), we get

$$\begin{aligned} \epsilon^{(l)} &= \Theta(1) \sum_{m=1}^l m^{-\alpha} + \epsilon^{(0)} \Theta(1) \\ &= \begin{cases} \Theta(l^{1-\alpha}) & \text{if } \alpha \in (0, 1) \\ \Theta(\log l) & \text{if } \alpha = 1 \\ \Theta(1) & \text{if } \alpha > 1. \end{cases} \end{aligned}$$

Now suppose  $\beta = 1$ . By [Lemma C.10](#), we get

$$\epsilon^{(l)} = \sum_{m=1}^l m^{-\alpha} \Theta(m^{-\delta} l^\delta) + \epsilon^{(0)} \Theta(l^\delta)$$

where the constants hidden inside the  $\Theta$  are the same in every term of the sum. If  $\alpha > 1 - \delta$ , then  $m^{-\delta-\alpha} = o(m^{-1})$ , so that  $\sum_{m=1}^l m^{-\delta-\alpha} = \Theta(1)$ , and

$$\begin{aligned} \epsilon^{(l)} &= \Theta(l^\delta) + \epsilon^{(0)} \Theta(l^\delta) \\ &= \Theta(l^\delta). \end{aligned}$$

On the other hand, if  $\alpha < 1 - \delta$ , then  $\sum_{m=1}^l m^{-\delta-\alpha} = \Theta(l^{1-\delta-\alpha})$ . So

$$\begin{aligned} \epsilon^{(l)} &= \Theta(l^{1-\alpha}) + \epsilon^{(0)} \Theta(l^\delta) \\ &= \Theta(l^{1-\alpha}). \end{aligned}$$

If  $\alpha = 1 - \delta$ , then  $\sum_{m=1}^l m^{-\delta-\alpha} = \Theta(\log l)$ . So

$$\begin{aligned} \epsilon^{(l)} &= \Theta(l^\delta \log l) + \epsilon^{(0)} \Theta(l^\delta) \\ &= \Theta(l^\delta \log l). \end{aligned}$$

Finally, if  $\beta \in (0, 1)$ , then

$$\epsilon^{(l)} = e^{\frac{\delta}{1-\beta} l^{1-\beta} + \Theta(l^{1-2\beta})} \sum_{m=1}^l m^{-\alpha} e^{\frac{-\delta}{1-\beta} m^{1-\beta} + \Theta(m^{1-2\beta})} + e^{\frac{\delta}{1-\beta} l^{1-\beta} + \Theta(l^{1-2\beta})}$$

The case of  $\delta = -1$  telescopes, so that the upper and lower constants hidden in  $\Theta$  can both be taken to be 1.  $\square$

**Lemma C.12.** Suppose for some  $\beta > 0$ , a sequence  $\epsilon^{(l)}$  satisfies

$$\epsilon^{(l)} = \epsilon^{(l-1)} (1 - \mu (\epsilon^{(l-1)})^\beta / l), \quad \epsilon^{(0)} \in (0, \frac{1}{\mu}).$$

Then  $\epsilon^{(l)} \sim (\beta \mu \log l)^{-1/\beta}$ .

*Proof.* Consider the differential equation

$$\dot{x}_\mu = -\mu x_\mu^{\beta+1}/t$$

for constant  $\mu$  has solution  $x_\mu = [\beta(\mu \log t + C)]^{-1/\beta}$  for some constant  $C$  determined by initial condition. Note that

$$-\mu x_\mu(t)^{\beta+1}/t \leq x_\mu(t+1) - x_\mu(t) \leq -\mu x_\mu(t+1)^{\beta+1}/(t+1) = -(1 - o(t^{-1}))\mu x_\mu(t)^{\beta+1}/t.$$

For any small enough  $\alpha > 0$ , we apply **Lemma C.9** with  $F(\epsilon, l) = \epsilon - \mu \epsilon^{\beta+1}/l$  (which is monotonic in  $\epsilon$  for small enough  $\epsilon$ ),  $c^{(l)} = x_\mu(l)$ , and  $b^{(l)} = x_{\mu-\alpha}(l)$  to obtain

$$x_{\mu-\alpha}(l) \leq \epsilon^{(l)} \leq x_\mu(l)$$

for large enough  $l$  and appropriately chosen initial conditions. This shows that  $\epsilon^{(l)} = \Theta(\log l^{-1/\beta})$ . Taking  $\alpha \rightarrow 0$ , we also obtain the leading coefficient  $\epsilon^{(l)} \sim [\beta \mu \log l]^{-1/\beta}$ . □

**Lemma C.13.** *Suppose a sequence  $u^{(l)}$  is governed by the equation*

$$u^{(l)} - u^{(l-1)} = A(u^{(l-1)} + B)^\alpha,$$

where  $\alpha \in [0, 1)$  and  $A > 0$ . Then  $u^{(l)} = K_1 l^{\frac{1}{1-\alpha}} - K_2 l^{\frac{\alpha}{1-\alpha}} \log l + o(l^{\frac{\alpha}{1-\alpha}} \log l)$ , where  $K_1 = [A(1-\alpha)]^{\frac{1}{1-\alpha}}$  and  $K_2 = \frac{1}{2} A^{\frac{1}{1-\alpha}} (1-\alpha)^{\frac{1}{1-\alpha}-1} \alpha$ .

*Proof. Leading term.* The differential equation

$$\dot{x}_{A,B} = A(x_{A,B} + B)^\alpha$$

has solution  $x_{A,B}(l) = [A(1-\alpha)(l+S)]^{\frac{1}{1-\alpha}} - B$  for some constant  $S$ . Since  $\dot{x}_{A,B}$  is monotonic, we have (writing  $x = x_{A,B}$  for brevity)

$$A(x_{A,B}(l) + B)^\alpha = \dot{x}_{A,B}(l) \leq x_{A,B}(l+1) - x_{A,B}(l) \leq \dot{x}_{A,B}(l+1) \leq (A + o(1))(x_{A,B}(l) + B)^\alpha$$

for large enough  $l$ . We apply **Lemma C.9** with  $F(x, l) = x + A(x+B)^\alpha$  (which is monotonic in  $x$  for large  $x$ ),  $c^{(l)} = x_{A,B}(l)$ , and  $b^{(l)} = x_{A-\epsilon, B}(l)$  to obtain

$$x_{A-\epsilon, B}(l) \leq u^{(l)} \leq x_{A, B}(l)$$

for large enough  $l$  and appropriate initial conditions. Therefore  $\lim u^{(l)}/l^{\frac{1}{1-\alpha}} \in [[(A-\epsilon)(1-\alpha)]^{\frac{1}{1-\alpha}}, [A(1-\alpha)]^{\frac{1}{1-\alpha}}]$ . Taking  $\epsilon \rightarrow 0$  gives the leading term.

**Subleading term.** Now let  $v^{(l)} := u^{(l)} - \aleph l^{\frac{1}{1-\alpha}}$ , where  $\aleph = [A(1-\alpha)]^{\frac{1}{1-\alpha}}$ . Then we have the recurrence

$$\begin{aligned} v^{(l+1)} + \aleph(l+1)^{\frac{1}{1-\alpha}} - v^{(l)} - \aleph l^{\frac{1}{1-\alpha}} &= A(v^{(l)} + \aleph l^{\frac{1}{1-\alpha}} + B)^\alpha \\ v^{(l+1)} - v^{(l)} + \aleph \left( \frac{1}{1-\alpha} l^{\frac{\alpha}{1-\alpha}} + \frac{1}{2} \left( \frac{1}{1-\alpha} \right) \left( \frac{\alpha}{1-\alpha} \right) l^{\frac{\alpha}{1-\alpha}-1} + \Theta(l^{\frac{\alpha}{1-\alpha}-2}) \right) \\ &= A[\aleph^\alpha l^{\frac{\alpha}{1-\alpha}} + \alpha(v^{(l)} + B)\aleph^{\alpha-1} l^{-1} + \Theta((v^{(l)} + B)l^{-1-\frac{1}{1-\alpha}})] \\ v^{(l+1)} - v^{(l)} &= \frac{\alpha}{1-\alpha} v^{(l)} l^{-1} - \frac{1}{2} \aleph \left( \frac{1}{1-\alpha} \right) \left( \frac{\alpha}{1-\alpha} \right) l^{\frac{\alpha}{1-\alpha}-1} + g(l) \end{aligned}$$

for some  $g(l) = O(l^{\frac{\alpha}{1-\alpha}-2} + l^{-1})$  and where, to get the last equation, we have used  $A\alpha^\alpha = \frac{1}{1-\alpha}\aleph$  to cancel the  $l^{\frac{\alpha}{1-\alpha}}$  term and simplified  $\alpha A \aleph^{\alpha-1} = \frac{\alpha}{1-\alpha}$ .

For any  $J > 0$ , the differential equation  $\dot{v}_J(l) = \frac{\alpha}{1-\alpha} v_J(l) l^{-1} - J l^{\frac{\alpha}{1-\alpha}-1}$  has solution  $v_J(l) = C[l(1-\alpha)]^{\frac{\alpha}{1-\alpha}} - J l^{\frac{\alpha}{1-\alpha}} \log l$ . Note that the functions  $F_J(z, n) = z + \frac{\alpha}{1-\alpha} z n^{-1} - J n^{\frac{\alpha}{1-\alpha}-1}$  and  $G_J(z, n) = F_J(z, n) + g(n)$  is monotonic in  $z$  (for positive  $n$ ). For large  $l$ , we also have  $\dot{v}_J(l)$  and  $F_J(v_J(l), l) = v_J(l) + \dot{v}_J(l)$  decreasing in  $l$ . Thus for any  $\epsilon > 0$  and  $l$  large enough

$$G_{J+\epsilon}(v_J(l), l) \leq F_{J+\epsilon/2}(v_J(l), l) \leq v_J(l) + \dot{v}_J(l+1) \leq v_J(l+1) \leq F_J(v_J(l), l) \leq G_{J-\epsilon}(v_J(l), l).$$

Now apply [Lemma C.9](#) with  $F = G_K$ ,  $a^{(l)} = v^{(l)}$ ,  $c^{(l)} = v_{K-\epsilon}$ ,  $b^{(l)} = v_{K+\epsilon}$  where  $K := \frac{1}{2} \aleph(\frac{1}{1-\alpha})(\frac{\alpha}{1-\alpha}) = \frac{1}{2} A^{\frac{1}{1-\alpha}} (1-\alpha)^{\frac{\alpha}{1-\alpha}-1} \alpha$ , with appropriately chosen initial conditions. This yields  $\lim_{l \rightarrow \infty} v^{(l)} / (l^{\frac{\alpha}{1-\alpha}} \log l) \in [-K - \epsilon, -K + \epsilon]$  for every  $\epsilon > 0$ , and there it must be equal to  $K$ . We have thus obtained the asymptotic expansion

$$u^{(l)} = [A(1-\alpha)l]^{\frac{1}{1-\alpha}} - \frac{1}{2} A^{\frac{1}{1-\alpha}} (1-\alpha)^{\frac{\alpha}{1-\alpha}-1} \alpha l^{\frac{\alpha}{1-\alpha}} \log l + o(l^{\frac{\alpha}{1-\alpha}} \log l).$$

□

**Lemma C.14.** *Suppose a sequence  $u^{(l)}$  is governed by the equation*

$$u^{(l)} - u^{(l-1)} = -A(u^{(l-1)} + B)^\alpha,$$

where  $\alpha > 1$  and  $A > 0$ . Then  $u^{(l)} \sim [A(\alpha - 1)l]^{\frac{1}{1-\alpha}}$ .

*Proof.* Similar to [Lemma C.13](#).

□

**Lemma C.15.** *Suppose a sequence  $u^{(l)}$  is governed by the equation*

$$u^{(l)} - u^{(l-1)} = A(u^{(l-1)} + B)^\alpha + C,$$

where  $\alpha \in (0, 1)$ . Then  $u^{(l)} = K_1 l^{\frac{1}{1-\alpha}} + R(l)$ , where the remainder  $R(l)$  is

$$R(l) \sim \begin{cases} -K_2 l^{\frac{\alpha}{1-\alpha}} \log l & \text{if } \alpha > \frac{1}{2} \\ (C - K_2) l \log l & \text{if } \alpha = \frac{1}{2} \text{ and } K_2 \neq C \\ \frac{C(1-\alpha)}{1-2\alpha} l & \text{if } \alpha < \frac{1}{2} \end{cases}$$

where  $K_1 = [A(1-\alpha)]^{\frac{1}{1-\alpha}}$ ,  $K_2 = \frac{1}{2} A^{\frac{1}{1-\alpha}} (1-\alpha)^{\frac{\alpha}{1-\alpha}-1} \alpha$  as in [Lemma C.13](#).

*Proof.*  $u$  is bounded below by the dynamics  $v^{(l)} - v^{(l-1)} = A(v^{(l-1)} + B)^\alpha$  and bounded above by the dynamics  $w^{(l)} - w^{(l-1)} = (A + o(1))(w^{(l-1)} + B)^\alpha$ . By [Lemma C.13](#), both  $v$  and  $w$  are asymptotic to  $u^{(l)} \sim [A(1-\alpha)l]^{\frac{1}{1-\alpha}}$ , which gives the result.

Now define  $v^{(l)} = u^{(l)} - [A(1-\alpha)l]^{\frac{1}{1-\alpha}}$ , and similar to the proof of [Lemma C.13](#), we find

$$v^{(l+1)} - v^{(l)} = \frac{\alpha}{1-\alpha} v^{(l)} l^{-1} - K l^{\frac{\alpha}{1-\alpha}-1} + C + g(l)$$

where  $K = \frac{1}{2} A^{\frac{1}{1-\alpha}} (1-\alpha)^{\frac{\alpha}{1-\alpha}-1} \alpha$  and  $g(l) = O(l^{\frac{\alpha}{1-\alpha}-2} + l^{-1})$ . If  $\frac{\alpha}{1-\alpha} > 1 \iff \alpha > \frac{1}{2}$ , then  $C + g(l) = o(l^{\frac{\alpha}{1-\alpha}-1})$  and we can proceed as in the proof of [Lemma C.13](#) to find  $v^{(l)} \sim K l^{\frac{\alpha}{1-\alpha}} \log l$ . If  $\frac{\alpha}{1-\alpha} = 1 \iff \alpha = 1$  and  $K \neq C$ , then  $v^{(l+1)} - v^{(l)} = \frac{\alpha}{1-\alpha} v^{(l)} l^{-1} - (K - C) l^{\frac{\alpha}{1-\alpha}-1} + g(l)$ , so that the technique used in [Lemma C.13](#) would obtain  $v^{(l)} \sim (K - C) l^{\frac{\alpha}{1-\alpha}} \log l = (K - C) l \log l$ . If  $\frac{\alpha}{1-\alpha} < 1 \iff \alpha < \frac{1}{2}$ , then  $v^{(l+1)} - v^{(l)} = \frac{\alpha}{1-\alpha} v^{(l)} l^{-1} + C + o(1)$ , then by using the differential equation  $\dot{v}_J(l) = \frac{\alpha}{1-\alpha} v_J(l) l^{-1} + J$  to approximate the difference equation solution and applying [Lemma C.9](#) as in the proof of [Lemma C.13](#), we obtain  $v^{(l)}(l) \sim \frac{C(1-\alpha)}{1-2\alpha} l$ . □

### C.3 Forward Dynamical Equations

Here we derive the recurrences governing the forward length and correlation quantities  $\mathbf{p}$ ,  $\mathbf{q}$ ,  $\boldsymbol{\lambda}$ ,  $\boldsymbol{\gamma}$ . We start with reduced residual networks.

**Lemma B.1.** *Suppose  $\phi$  is antisymmetric. Then in an RRN,  $\mathbf{p}$  and  $\mathbf{q}$  satisfy the recurrence*

$$\begin{aligned} \mathbf{q} &= \sigma_w^2 \underline{\mathbf{p}} + \sigma_b^2 \\ \mathbf{p} &= \mathbb{V} \phi(\mathbf{q}) + \underline{\mathbf{p}}. \end{aligned}$$

*Proof.* We have

$$\begin{aligned}\mathbf{q} &= \langle h_j^2 \rangle = \left\langle \sum_i (w_{ji} \underline{x}_i + b_j)^2 \right\rangle \\ &= \langle b_j^2 \rangle + \sum_i \langle w_{ji}^2 \underline{x}_i^2 \rangle + 2 \sum_i \langle w_{ji} \underline{x}_i b_j \rangle + 2 \sum_{j \neq l} \langle w_{ji} w_{li} \underline{x}_i^2 \rangle\end{aligned}$$

But  $w_{ji}$ ,  $w_{li}$ ,  $\underline{x}$ , and  $b_j$  form an independency, so the last two sums are 0, and the terms in the first sum split multiplicatively. Therefore

$$\begin{aligned}\mathbf{q} &= \sigma_b^2 + \sum_i \langle w_{ji}^2 \rangle \langle \underline{x}_i^2 \rangle \\ &= \sigma_b^2 + N \cdot \frac{\sigma_w^2}{N} \underline{\mathbf{p}} \\ &= \sigma_b^2 + \sigma_w^2 \underline{\mathbf{p}}.\end{aligned}$$

For the recurrence of  $\underline{\mathbf{p}}$ , we have

$$\begin{aligned}\mathbf{p} &= \langle x_i^2 \rangle = \langle (\phi(h_i) + \underline{x}_i)^2 \rangle \\ &= \langle \phi(h_i)^2 \rangle + \langle \underline{x}_i^2 \rangle + 2 \langle \phi(h_i) \underline{x}_i \rangle\end{aligned}$$

As  $N \rightarrow \infty$ , the coefficient  $w_{ji}$  of  $\underline{x}_i$  in  $h_i$  has vanishing covariance, so  $h_i$  and  $\underline{x}_i$  become independent. Therefore  $\langle \phi(h_i) \underline{x}_i \rangle = \langle \phi(h_i) \rangle \langle \underline{x}_i \rangle$ . Because  $h_i$  is the sum of a large number of independent random variables, by CLT,  $h_i$  is a Gaussian with mean  $\sum_i \langle w_{ji} \rangle \langle \underline{x}_i \rangle + \langle b_j \rangle = 0$  since  $\langle w_{ji} \rangle = \langle b_j \rangle = 0$ . Our antisymmetry assumption on  $\phi$  then implies  $\langle \phi(h_i) \rangle = 0$ . Therefore,

$$\begin{aligned}\mathbf{p} &= \langle \phi(h_i)^2 \rangle + \langle \underline{x}_i^2 \rangle \\ &= \mathbf{V}\phi(\mathbf{q}) + \underline{\mathbf{p}}\end{aligned}$$

as desired.  $\square$

**Theorem B.3.** Suppose  $\phi$  is antisymmetric. Then in an RRN,  $\lambda$  and  $\gamma$  satisfy the recurrence

$$\begin{aligned}\lambda &= \sigma_w^2 \underline{\gamma} + \sigma_b^2 \\ \gamma &= \mathbf{W}\phi(\mathbf{q}, \lambda) + \underline{\gamma}.\end{aligned}$$

*Proof.* Similar to [Lemma B.1](#).  $\square$

Now, for the full residual networks, the proofs are similar, but we no longer need to assume that  $\phi$  is antisymmetric because of the randomization via the extra sets of weights.

**Theorem B.8.** For any nonlinearity  $\phi$  in an FRN,

$$\begin{aligned}\mathbf{q} &= \sigma_w^2 \underline{\mathbf{p}} + \sigma_b^2 \\ \mathbf{p} &= \sigma_v^2 \mathbf{V}\phi(\mathbf{q}) + \sigma_a^2 + \underline{\mathbf{p}}\end{aligned}$$

*Proof.*

$$\begin{aligned}\mathbf{q} &= \langle h_j^2 \rangle = \langle (w_j^i \underline{x}_i + b_j)^2 \rangle = \langle (w_j^i \underline{x}_i)^2 \rangle + \langle b_j^2 \rangle \\ &= \sigma_w^2 \langle \underline{x}_i^2 \rangle + \sigma_b^2 \\ &= \sigma_w^2 \underline{\mathbf{p}} + \sigma_b^2 \\ \mathbf{p} &= \langle x_i^2 \rangle = \langle (v_i^j \phi(h_j) + \underline{x}_i + a_i)^2 \rangle \\ &= \sigma_v^2 \langle \phi(h_i)^2 \rangle + \langle \underline{x}_i^2 \rangle + \sigma_a^2 \\ &= \sigma_v^2 \mathbf{V}\phi(\mathbf{q}) + \sigma_a^2 + \underline{\mathbf{p}}\end{aligned}$$

where in the third equality for  $\underline{\mathbf{p}}$ , we are now using the independence of  $v_i^j$  from all other variables to cancel out the terms, whereas before we had to rely on  $\phi$  being antisymmetric.  $\square$



**Theorem B.10.** For any nonlinearity  $\phi$ , in an FRN

$$\begin{aligned}\lambda &= \sigma_w^2 \underline{\gamma} + \sigma_b^2 \\ \gamma &= \sigma_v^2 \mathbf{W} \phi(\mathbf{q}, \lambda) + \sigma_a^2 + \underline{\gamma}\end{aligned}$$

*Proof.* Similar to **Thm B.8**. □

#### C.4 Backward Dynamical Equations

Here we derive the recurrences governing the gradient quantities  $\chi$  and  $\chi_\bullet$  for different  $\bullet$ , all under the gradient independence assumption. Write  $\beta_i^{(l)} = \frac{\partial E}{\partial x_i^{(l)}}$  for a cost function  $E$ .

**Theorem B.5.** For any nonlinearity  $\phi$  in an RRN, under assumptions **Axiom 3.1** and **Axiom 3.2**, whenever  $\dot{\phi}^2(\zeta)$  has finite variance for Gaussian variable  $\zeta$ ,

$$\underline{\chi} = (\sigma_w^2 \mathbf{V} \dot{\phi}(\mathbf{q}) + 1) \chi, \quad \chi_b = \chi \mathbf{V} \dot{\phi}(\mathbf{q}), \quad \chi_w = \chi \mathbf{V} \dot{\phi}(\mathbf{q}) \underline{\mathbf{p}}.$$

*Proof.* For a reduced residual network, we have the following derivative computation:

$$\frac{\partial x_i}{\partial \underline{x}_j} = \delta_{ji} + \dot{\phi}(h_i) \frac{\partial h_i}{\partial \underline{x}_j}, \quad \frac{\partial x_i}{\partial h_j} = \delta_{ji} \dot{\phi}(h_j), \quad \frac{\partial h_i}{\partial \underline{x}_j} = w_{ij}, \quad \frac{\partial h_i}{\partial w_{ij}} = \underline{x}_j, \quad \frac{\partial h_i}{\partial b_j} = \delta_{ij}.$$

Then

$$\begin{aligned}\underline{\beta}_j &= \beta_j + \sum_i \beta_i \dot{\phi}(h_i) \frac{\partial h_i}{\partial \underline{x}_j} \\ &= \beta_j + \sum_i \beta_i \dot{\phi}(h_i) w_{ij} \\ \langle \underline{\beta}_j^2 \rangle &= \langle [\beta_j + \sum_i \beta_i \dot{\phi}(h_i) w_{ij}]^2 \rangle \\ &= \langle \beta_j^2 \rangle + \sum_i \langle \beta_i^2 \dot{\phi}^2(h_i) (w_{ij})^2 \rangle \\ &\quad + 2 \sum_{i < k} \langle \beta_i \beta_k \dot{\phi}(h_i) w_{ij} \dot{\phi}(h_k) w_{kj} \rangle + 2 \sum_i \langle \beta_j \beta_i \dot{\phi}(h_i) w_{ij} \rangle\end{aligned}$$

The last two terms of the above vanish as  $w_{ij}$  is independent from  $w_{kj}$ ,  $h_i, h_k$  and  $\beta_i, \beta_j, \beta_k$  by **Axiom 3.2**, and  $\langle w_{ij} \rangle = 0$ .

Therefore, applying **Axiom 3.1**,

$$\begin{aligned}\langle \underline{\beta}_j^2 \rangle &= \sigma_w^2 \langle \beta_j^2 \rangle \langle \dot{\phi}^2(h_i) \rangle + \langle \beta_j^2 \rangle \\ &= (\sigma_w^2 \mathbf{V} \dot{\phi}(\mathbf{q}) + 1) \langle \beta_j^2 \rangle\end{aligned}$$

We similarly have

$$\begin{aligned}\frac{\partial E}{\partial b_j} &= \sum_i \frac{\partial E}{\partial x_i} \frac{\partial x_i}{\partial h_j} = \beta_j \dot{\phi}(h_j), & \text{since } \frac{\partial x_i}{\partial h_j} &= \delta_{ji} \dot{\phi}(h_j) \\ \left\langle \left( \frac{\partial E}{\partial b_j} \right)^2 \right\rangle &= \langle \beta_j^2 \dot{\phi}^2(h_j) \rangle = \langle \beta_j^2 \rangle \mathbf{V} \dot{\phi}(\mathbf{q}), & \text{by } \mathbf{Axiom 3.2(b)}; \\ \frac{\partial E}{\partial w_{ji}} &= \sum_i \frac{\partial E}{\partial x_i} \frac{\partial x_i}{\partial h_j} \frac{\partial h_j}{\partial w_{ji}} = \beta_j \dot{\phi}(h_j) \underline{x}_i, & \text{since } \frac{\partial x_i}{\partial h_j} &= \delta_{ji} \dot{\phi}(h_j) \\ \left\langle \left( \frac{\partial E}{\partial w_{ji}} \right)^2 \right\rangle &= \langle \beta_j^2 \dot{\phi}^2(h_j) \underline{x}_i^2 \rangle = \langle \beta_j^2 \rangle \mathbf{V} \dot{\phi}(\mathbf{q}) \underline{\mathbf{p}}, & \text{by } \mathbf{Axiom 3.2(b)}\end{aligned}$$

In the last equation we have also used the fact that as  $N \rightarrow \infty$ ,  $h_j$  and  $x_i$  become independent (they are jointly Gaussian and their correlation  $\langle w_{ji}^2 \rangle$  goes to 0 with  $N$ ). □

**Theorem B.12.** For any nonlinearity  $\phi$  in an FRN, under assumptions [Axiom 3.1](#) and [Axiom 3.2](#), whenever  $\dot{\phi}(\zeta)^2$  has finite variance for Gaussian variable  $\zeta$ ,

$$\begin{aligned}\underline{\chi} &= (\sigma_v^2 \sigma_w^2 \mathbf{V} \dot{\phi}(\mathbf{q}) + 1) \underline{\chi}, & \chi_b &= \sigma_v^2 \chi \mathbf{V} \dot{\phi}(\mathbf{q}), \\ \chi_w &= \sigma_v^2 \chi \mathbf{V} \dot{\phi}(\mathbf{q}) \underline{\mathbf{p}}, & \chi_v &= \chi \mathbf{V} \dot{\phi}(\mathbf{q}), & \chi_a &= \chi\end{aligned}$$

*Proof.* For the full residual network, we have the following derivative computations:

$$\begin{aligned}\frac{\partial x_i}{\partial \underline{x}_j} &= \delta_{ji} + \sum_k v_{ik} \dot{\phi}(h_k) \frac{\partial h_k}{\partial \underline{x}_j}, & \frac{\partial x_i}{\partial h_j} &= v_{ij} \dot{\phi}(h_j), & \frac{\partial h_i}{\partial \underline{x}_j} &= w_{ij}, & \frac{\partial h_i}{\partial w_{ij}} &= \underline{x}_j, & \frac{\partial h_i}{\partial b_i} &= 1, \\ & & \frac{\partial x_i}{\partial v_{ik}} &= \phi(h_k), & \frac{\partial x_i}{\partial a_i} &= 1.\end{aligned}$$

Again let  $\beta_j = \frac{\partial E}{\partial x_j}$ . Then

$$\begin{aligned}\underline{\beta}_j &= \sum_i \beta_i (\delta_{ji} + \sum_k v_{ik} \dot{\phi}(h_k) \frac{\partial h_k}{\partial \underline{x}_j}) \\ &= \sum_i \beta_i (\delta_{ji} + \sum_k v_{ik} \dot{\phi}(h_k) w_{kj})\end{aligned}$$

Thus,

$$\begin{aligned}\langle \underline{\beta}_j^2 \rangle &= \langle [\sum_i \beta_i (\delta_{ji} + \sum_k v_{ik} \dot{\phi}(h_k) w_{kj})]^2 \rangle \\ &= \langle \beta_j^2 \rangle + \sum_{i,k} \langle v_{ik}^2 \rangle \langle w_{kj}^2 \rangle \mathbf{V} \dot{\phi}(\mathbf{q}) \langle \beta_i^2 \rangle \\ &= \langle \beta_j^2 \rangle (1 + \sigma_v^2 \sigma_w^2 \mathbf{V} \dot{\phi}(\mathbf{q}))\end{aligned}$$

where in the second equality we applied the independence argument as in the proof of [Thm B.5](#), leveraging [Axiom 3.2](#), and in the third equality we used [Axiom 3.1](#) to get  $\langle \beta_i^2 \rangle = \langle \beta_j^2 \rangle$ .

The other computations are similar to the proof of [Thm B.12](#). □

## C.5 Tanh: Reduced Residual Network

### C.5.1 Forward Dynamics

**Theorem B.2.** Suppose  $\phi$  is tanh-like. Assume RRN architecture.

- If  $\sigma_w = 0$ , then  $\mathbf{p}^{(l)} = l \mathbf{V} \phi(\sigma_b^2) + \mathbf{p}^{(0)}$  and  $\mathbf{q}^{(l)} = \sigma_b^2$ .
- If  $\sigma_w > 0$ ,  $\lim_{l \rightarrow \infty} \mathbf{p}^{(l)}/l = 1$  and  $\lim_{l \rightarrow \infty} \mathbf{q}^{(l)}/(\sigma_w^2 l) = 1$ . If  $\phi = \tanh$ , then we can obtain more terms of the asymptotic expansions:

$$\begin{aligned}\mathbf{p}^{(l)} &= l - 2C \sigma_w^{-1} l^{1/2} - C^2 \sigma_w^{-2} \log l + O(1) \\ \mathbf{q}^{(l)} &= \sigma_w^2 l - 2C \sigma_w l^{1/2} - C^2 \log l + O(1)\end{aligned}$$

as  $l \rightarrow \infty$ , where  $C = \sqrt{2/\pi}$ .

*Proof.* The case with  $\sigma_w = 0$  is trivial. We assume  $\sigma_w > 0$  from here on.

$\mathbf{p}$  and  $\mathbf{q}$  are asymptotically linear with  $l$ . We first show that, for any  $\omega < 1$ ,

$$l + \mathbf{p}^{(0)} \geq \mathbf{p}^{(l)} \geq \omega l$$

and

$$\sigma_w^2 (l + \mathbf{p}^{(0)}) + \sigma_b^2 \geq \mathbf{q}^{(l)} \geq \sigma_w^2 \omega (l - 1) + \sigma_b^2,$$

so that  $\mathbf{p}^{(l)} \sim l$  and  $\mathbf{q}^{(l)} \sim \sigma_w^2 l$ .

The upper bounds are trivial, given  $V\phi(\mathbf{q}) \leq 1$  for any  $\mathbf{q}$ . We show the lower bounds for any  $\omega < 1$ .

For any  $\epsilon > 0$ , define  $\aleph_\epsilon$  by  $\phi^2(\aleph_\epsilon) = \exp(-\epsilon)$ . Then

$$\begin{aligned} V\phi(\mathbf{q}) &\geq \exp(-\epsilon) \Pr[z \notin [-\aleph_\epsilon, \aleph_\epsilon] : z \sim \mathcal{N}(0, \mathbf{q})] \\ &\geq \exp(-\epsilon) \left(1 - \frac{2\aleph_\epsilon}{\sqrt{2\pi\mathbf{q}}}\right) \end{aligned}$$

where the second inequality follows from an overestimate of the  $\Pr[z \in [-\aleph_\epsilon, \aleph_\epsilon]]$  via the mode of  $\mathcal{N}(0, \mathbf{q})$ .

For any  $\mathbf{q} \geq \mathbf{q}^{(0)}$ ,  $V\phi(\mathbf{q})$  is then lower bounded by

$$\phi^2\left(\sqrt{\mathbf{q}^{(0)}}\right) \left(1 - \frac{2\sqrt{\mathbf{q}^{(0)}}}{\sqrt{2\pi\mathbf{q}^{(0)}}}\right) = \phi^2\left(\sqrt{\mathbf{q}^{(0)}}\right) \left(1 - \sqrt{\frac{2}{\pi}}\right) > 0.$$

Thus  $\mathbf{p}^{(l)}$  and  $\mathbf{q}^{(l)}$  are unbounded with  $l$ .

Furthermore, as  $\mathbf{q} \rightarrow \infty$ , the lower bound  $\exp(-\epsilon) \left(1 - \frac{2\aleph_\epsilon}{\sqrt{2\pi\mathbf{q}}}\right)$  goes to  $\exp(-\epsilon)$ , for any  $\epsilon$ . Therefore, for any  $\omega < 1$ ,  $\mathbf{p}^{(l)} \geq \omega l$  and  $\mathbf{q}^{(l)} \geq \sigma_w^2 \omega (l-1) + \sigma_b^2$ .

**Asymptotic expansion.** Now we repeat the following to get each successive asymptotic term of  $\mathbf{p}^{(l)}$  and  $\mathbf{q}^{(l)}$ : We plug in the current asymptotic form of  $\mathbf{q}^{(l)}$  into  $V \tanh(\mathbf{q}) = 1 - C\mathbf{q}^{-1/2} + \Theta(\mathbf{q}^{-3/2})$  ([Lemma C.5](#)), where  $C = \sqrt{2/\pi}$ . Next we take the sum  $\mathbf{q}^{(l)} = \sum_{r=1}^l V \tanh(\mathbf{q}^{(r)})$ , which yields one more term in the asymptotic expansion of  $\mathbf{p}$  than the last round. We then repeat until we get only constant terms.

The following exhibits a trace of this procedure, where in the summation step for  $\mathbf{q}^{(l)}$ , we implicitly apply

$$\begin{aligned} \mathbf{q} &= \sigma_w^2 l + o(l) = \sigma_w^2 l(1 + o(1)) \\ \mathbf{q}^{-1/2} &= \sigma_w^{-1} l^{-1/2} (1 + o(1)) = \sigma_w^{-1} l^{-1/2} + o(l^{-1/2}) \\ \mathbf{p} &= \sum_{r=1}^l 1 - C(\mathbf{q}^{(r)})^{-1/2} + \Theta((\mathbf{q}^{(r)})^{-3/2}) \\ &= \sum_{r=1}^l 1 - C(\sigma_w^{-1} r^{-1/2} + o(r^{-1/2})) + \Theta(r^{-3/2}) \\ &= l - 2C\sigma_w^{-1} l^{1/2} + o(l^{1/2}) \\ \mathbf{q} &= \sigma_w^2 l - 2C\sigma_w l^{1/2} + o(l^{1/2}) = \sigma_w^2 l(1 - 2C\sigma_w^{-1} l^{-1/2} + o(l^{-1/2})) \\ \mathbf{q}^{-1/2} &= \sigma_w^{-1} l^{-1/2} (1 + C\sigma_w^{-1} l^{-1/2} + o(l^{-1/2})) = \sigma_w^{-1} l^{-1/2} + C\sigma_w^{-2} l^{-1} + o(l^{-1}) \\ \mathbf{p} &= \sum_{r=1}^l 1 - C(\sigma_w^{-1} l^{-1/2} + C\sigma_w^{-2} l^{-1} + o(l^{-1})) + \Theta(l^{-3/2}) \\ &= l - 2C\sigma_w^{-1} l^{1/2} - C^2 \sigma_w^{-2} \log l + o(\log l) \\ \mathbf{q} &= \sigma_w^2 l \left(1 - 2C\sigma_w^{-1} l^{-1/2} - C^2 \sigma_w^{-2} \frac{\log l}{l} + o\left(\frac{\log l}{l}\right)\right) \\ \mathbf{q}^{-1/2} &= \sigma_w^{-1} l^{-1/2} \left(1 + C\sigma_w^{-1} l^{-1/2} + \frac{1}{2} C^2 \sigma_w^{-2} \frac{\log l}{l} + o\left(\frac{\log l}{l}\right)\right) \\ \mathbf{p} &= \sum_{r=1}^l 1 - C(\sigma_w^{-1} r^{-1/2} + C\sigma_w^{-2} r^{-1} + \frac{1}{2} C^2 \sigma_w^{-3} \frac{\log r}{r^{3/2}} + o\left(\frac{\log r}{r^{3/2}}\right)) + \Theta(r^{-3/2}) \\ &= l - 2C\sigma_w^{-1} l^{1/2} - C^2 \sigma_w^{-2} \log l + O(1) \end{aligned}$$

which is what we want.  $\square$

**Lemma C.16.** Let  $\phi$  is antisymmetric. Then for  $\tau \in [0, \pi/2]$ ,

$$\begin{aligned} \text{W}\phi(q, q \cos \tau) &= \lim_{t \rightarrow \tau} \frac{1}{\pi \sin t} \int_{w' \geq |w|} dw dw' \Upsilon(w, w'; \tau) \phi\left(\frac{\sqrt{q}}{\sqrt{2}}(w + w')\right) \phi\left(\frac{\sqrt{q}}{\sqrt{2}}(w' - w)\right) \\ &= \frac{1}{\pi} \int_0^\infty r dr e^{-r^2/2} \int_0^\pi d\theta \Sigma(\sqrt{qr}, \theta; \tau) \\ &= \frac{1}{\pi} \int_0^\infty s ds q^{-1} e^{-s^2 q^{-1}/2} \int_0^\pi d\theta \Sigma(s, \theta; \tau) \\ &= \frac{1}{\pi} \int_0^\pi d\theta \int_0^\infty ds e^{-s^2 q^{-1}/2} \frac{\partial}{\partial s} \Sigma(s, \theta; \tau) \end{aligned}$$

where  $\Upsilon(w, w'; \tau) := e^{-\frac{1}{2}(\frac{w^2}{1-c} + \frac{(w')^2}{1+c})} - e^{-\frac{1}{2}(\frac{(w')^2}{1-c} + \frac{w^2}{1+c})}$  with  $c = \cos \tau$ , and  $\Sigma(s, \theta; \tau) := \phi(s \sin \theta) \phi(s \sin(\theta - \tau))$ .

Of course, in the above lemma, the limit in the first equation is only necessary when  $\tau = 0$  or  $\tau = \pi/2$ .

*Proof.* Let  $c := \cos \tau$  and

$$\Gamma := \text{W}\phi(q, cq) = \frac{1}{2\pi q \sqrt{1-c^2}} \int d\mathbf{z} \exp(-\mathbf{z}^T \Sigma^{-1} \mathbf{z}/2) \phi(z) \phi(z'),$$

where  $\Sigma = \begin{pmatrix} q & cq \\ cq & q \end{pmatrix}$ .

Our proof will have two portions: Symmetrization of the  $\Gamma$  integral and trigonometric change of variables for evaluation.

**Symmetrization.**  $\Sigma$  is diagonalized by  $\Omega = \frac{1}{\sqrt{2q}} \begin{pmatrix} -1 & 1 \\ 1 & 1 \end{pmatrix}$ ,

$$\Sigma = \Omega^T \text{Diag}(1-c, 1+c) \Omega.$$

By a change of variable  $\mathbf{w} = \Omega \mathbf{z}$ , so that  $d\mathbf{w} = q^{-1} d\mathbf{z}$ , we have

$$\begin{aligned} \Gamma &= \frac{1}{2\pi \sqrt{1-c^2}} \int d\mathbf{w} \exp(-\mathbf{w}^T \text{Diag}(1-c, 1+c)^{-1} \mathbf{w}/2) \phi\left(\frac{\sqrt{q}}{\sqrt{2}}(w' - w)\right) \phi\left(\frac{\sqrt{q}}{\sqrt{2}}(w + w')\right) \\ &= \frac{1}{2\pi \sqrt{1-c^2}} \int dw dw' e^{-\frac{1}{2}(\frac{w^2}{1-c} + \frac{(w')^2}{1+c})} \phi\left(\frac{\sqrt{q}}{\sqrt{2}}(w' - w)\right) \phi\left(\frac{\sqrt{q}}{\sqrt{2}}(w + w')\right) \end{aligned}$$

By a change of variable swapping  $w$  with  $w'$ , we get

$$\Gamma = -\frac{1}{2\pi \sqrt{1-c^2}} \int dw dw' e^{-\frac{1}{2}(\frac{(w')^2}{1-c} + \frac{w^2}{1+c})} \phi\left(\frac{\sqrt{q}}{\sqrt{2}}(w + w')\right) \phi\left(\frac{\sqrt{q}}{\sqrt{2}}(w' - w)\right)$$

Thus

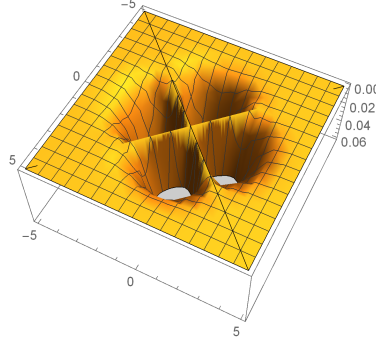
$$2\Gamma = \frac{1}{2\pi \sqrt{1-c^2}} \int dw dw' \Upsilon(w, w'; \tau) \phi\left(\frac{\sqrt{q}}{\sqrt{2}}(w + w')\right) \phi\left(\frac{\sqrt{q}}{\sqrt{2}}(w' - w)\right)$$

where

$$\Upsilon(w, w'; \tau) = e^{-\frac{1}{2}(\frac{w^2}{1-c} + \frac{(w')^2}{1+c})} - e^{-\frac{1}{2}(\frac{(w')^2}{1-c} + \frac{w^2}{1+c})}.$$

Note that, by the antisymmetry of  $\phi$ , the integrand  $K := \Upsilon(w, w'; \tau) \phi(\dots) \phi(\dots)$  above has the symmetries  $K(w, w') = K(w', w) = K(w, -w')$ , and is everywhere nonnegative. **Fig. C.10** displays a contour plot of  $K$  for typical values of  $q$  and  $c$ . So

$$\Gamma = \frac{1}{\pi \sqrt{1-c^2}} \int_{w' \geq |w|} dw dw' K(w, w').$$



**Figure C.10:** The integrand of  $\Gamma$  after symmetrization. Here  $c = .2$  and  $q = 100$  and  $\phi = \tanh$ .

This gives the first equation in the lemma.

**Polar Coordinates.** Let  $\frac{w}{\sqrt{1-c}} = r \cos \theta$ ,  $\frac{w'}{\sqrt{1+c}} = r \sin \theta$ , so that

$$\begin{aligned} w &= r \cos \theta \sqrt{1-c} = \sqrt{2} r \cos \theta \sin \frac{\tau}{2} \\ w' &= r \sin \theta \sqrt{1+c} = \sqrt{2} r \sin \theta \cos \frac{\tau}{2} \\ dw dw' &= \sqrt{1-c^2} r dr d\theta = (\sin^2 \tau) r dr d\theta. \end{aligned}$$

Then

$$\begin{aligned} \mathbf{A} &:= \int_{w' \geq |w|} e^{-\left(\frac{w^2}{1-c} + \frac{(w')^2}{1+c}\right)/2} \phi(\sqrt{q/2}(w+w')) \phi(\sqrt{q/2}(w'-w)) dw dw' \\ &= \sin^2 \tau \int_0^\infty r dr e^{-r^2/2} \int_{\tau/2}^{\pi-\tau/2} d\theta \phi(\sqrt{qr} \sin(\theta + \tau/2)) \phi(\sqrt{qr} \sin(\theta - \tau/2)). \end{aligned}$$

Similarly, let  $\frac{w}{\sqrt{1+c}} = r \cos \theta$ ,  $\frac{w'}{\sqrt{1-c}} = r \sin \theta$ , so that

$$\begin{aligned} w &= r \cos \theta \sqrt{1+c} = \sqrt{2} r \cos \theta \cos \frac{\tau}{2} \\ w' &= r \sin \theta \sqrt{1-c} = \sqrt{2} r \sin \theta \sin \frac{\tau}{2} \\ dw dw' &= \sqrt{1-c^2} r dr d\theta = (\sin^2 \tau) r dr d\theta, \end{aligned}$$

and

$$\begin{aligned} \mathbf{B} &= \int_{w' \geq |w|} e^{-\left(\frac{w^2}{1+c} + \frac{(w')^2}{1-c}\right)/2} \phi(\sqrt{q/2}(w+w')) \phi(\sqrt{q/2}(w'-w)) dw dw' \\ &= -\sin^2 \tau \int_0^\infty r dr e^{-r^2/2} \int_{\pi/2-\tau/2}^{\pi/2+\tau/2} d\theta \phi(\sqrt{qr} \cos(\theta + \tau/2)) \phi(\sqrt{qr} \cos(\theta - \tau/2)) \\ &= -\sin^2 \tau \int_0^\infty r dr e^{-r^2/2} \int_{-\tau/2}^{\tau/2} d\theta \phi(\sqrt{qr} \sin(\theta + \tau/2)) \phi(\sqrt{qr} \sin(\theta - \tau/2)). \end{aligned}$$

Thus

$$\begin{aligned} \Gamma &= \frac{1}{\pi \sqrt{1-c^2}} (\mathbf{A} - \mathbf{B}) \\ &= \frac{1}{\pi} \int_0^\infty r dr e^{-r^2/2} \int_{-\tau/2}^{\pi-\tau/2} d\theta \phi(\sqrt{qr} \sin(\theta + \tau/2)) \phi(\sqrt{qr} \sin(\theta - \tau/2)) \\ &= \frac{1}{\pi} \int_0^\infty r dr e^{-r^2/2} \int_0^\pi d\theta \phi(\sqrt{qr} \sin(\theta)) \phi(\sqrt{qr} \sin(\theta - \tau)). \end{aligned}$$

This gives the second equation in the lemma, and a change of variables  $s = \sqrt{q}r$  gives the third.

For the fourth equality, we start from the third equality, and apply integration by parts:

$$\begin{aligned}
& \frac{1}{\pi} \int_0^\infty s \, ds q^{-1} e^{-s^2 q^{-1}/2} \int_0^\pi d\theta \Sigma(s, \theta; \tau) \\
&= \frac{1}{\pi} \int_0^\pi d\theta \int_0^\infty ds s q^{-1} e^{-s^2 q^{-1}/2} \Sigma(s, \theta; \tau) \\
&= \frac{1}{\pi} \int_0^\pi d\theta \left( -e^{-s^2 q^{-1}/2} \Sigma(s, \theta; \tau) \Big|_{s=0}^\infty + \int_0^\infty ds e^{-s^2 q^{-1}/2} \frac{\partial}{\partial s} \Sigma(s, \theta; \tau) \right) \\
&= \frac{1}{\pi} \int_0^\pi d\theta \int_0^\infty ds e^{-s^2 q^{-1}/2} \frac{\partial}{\partial s} \Sigma(s, \theta; \tau).
\end{aligned}$$

where the last equality follows because  $\Sigma(0, \theta; \tau) = 0$  and  $e^{-s^2 q^{-1}/2} \rightarrow 0$  as  $s \rightarrow \infty$ .  $\square$

In the following lemmas, the “2” is not important, and can be any arbitrary finite or infinite value.

**Lemma C.17.** *Suppose a function  $f : (0, 2) \rightarrow \mathbb{R}$  is  $C^k$  on  $(0, 2)$ . If  $\lim_{x \downarrow 0} f^{(i)}(x)$  exists and is finite for every  $i \in [0, k]$ , then  $f$  can be extended to  $[0, 2)$  such that one sided  $i$ th derivatives exist at 0 for all  $i \in [0, k]$ .*

*Proof.* Consider  $\overline{f^{(i)}}(0) := f^{(i)}(1) - \int_0^1 f^{(i+1)}(x) \, dx$  for  $i \in [0, k-1]$ , which naturally is also equal to  $f^{(i)}(\epsilon) - \int_0^\epsilon f^{(i+1)}(x) \, dx$  for any  $\epsilon > 0$ . Certainly  $f^{(i)}(x) \rightarrow \overline{f^{(i)}}(0)$  as  $x \rightarrow 0$  if this limit exists — and by assumption it does, for  $0 \leq i \leq k-1$ . Therefore, we can define the extension of  $f^{(i)}$  to  $x = 0$  to be  $\overline{f^{(i)}}(0) := \overline{f^{(i)}}(0)$ . But we need to check that for  $i \in [0, k-1]$ .

$$\lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} (f^{(i)}(\epsilon) - f^{(i)}(0)) = f^{(i+1)}(0)$$

so that all one sided  $i$ th derivatives exist. But

$$\begin{aligned}
\frac{1}{\epsilon} (f^{(i)}(\epsilon) - f^{(i)}(0)) &= \frac{1}{\epsilon} \int_0^\epsilon f^{(i+1)}(x) \, dx \\
&= f^{(i+1)}(0) + \int_0^1 (f^{(i+1)}(x) - f^{(i+1)}(0)) \mathbf{I}(x \in [0, \epsilon]) \, dx
\end{aligned}$$

Since  $\lim_{x \downarrow 0} f^{(i+1)}(x) = f^{(i+1)}(0)$ ,  $f^{(i+1)}(x) - f^{(i+1)}(0)$  is bounded for small  $x$ , and by dominated convergence,  $\int_0^1 (f^{(i+1)}(x) - f^{(i+1)}(0)) \mathbf{I}(x \in [0, \epsilon]) \, dx \rightarrow \int_0^1 0 \, dx = 0$  as  $\epsilon \rightarrow 0$ . Thus

$$\lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} (f^{(i)}(\epsilon) - f^{(i)}(0)) = f^{(i+1)}(0)$$

as desired.  $\square$

**Lemma C.18.** *If  $f : [0, 2) \rightarrow \mathbb{R}$  is  $C^k$  on  $(0, 2)$  and has one sided derivatives at 0 up to order  $k$ , then*

$$f(\epsilon) = f(0) + \epsilon f^{(1)}(0) + \dots + \frac{\epsilon^{i-1}}{(i-1)!} f^{(i-1)}(0) + O(\epsilon^i)$$

for any  $i \leq k$ .

*Proof.* We have

$$\begin{aligned}
f(\epsilon) &= f(0) + \int_0^\epsilon f^{(1)}(x) dx \\
&= f(0) + \epsilon f^{(1)}(0) + \int_0^\epsilon f^{(1)}(x) - f^{(1)}(0) dx \\
&= f(0) + \epsilon f^{(1)}(0) + \int_0^\epsilon \int_0^{x_0} f^{(2)}(x_2) dx_2 dx_1 \\
&= f(0) + \epsilon f^{(1)}(0) + \frac{\epsilon^2}{2} f^{(2)}(0) + \int_0^\epsilon \int_0^{x_1} f^{(2)}(x_2) - f^{(2)}(0) dx_2 dx_1 \\
&\vdots \\
f(\epsilon) &= f(0) + \epsilon f^{(1)}(0) + \dots + \frac{\epsilon^{i-1}}{(i-1)!} f^{(i-1)}(0) + \int_0^\epsilon dx_1 \int_0^{x_1} dx_2 \dots \int_0^{x_{i-1}} dx_i f^{(i)}(x_i)
\end{aligned}$$

for any  $i \leq k$ . It suffices then to bound the size of the integral. Since  $f^{(i)}(x) \rightarrow f^{(i)}(0)$  as  $x \downarrow 0$  by assumption,  $|f^{(i)}(x_i)|$  is bounded by some constant  $C$  on the integration region  $\mathbb{A} := \{(x_1, \dots, x_i) : \epsilon \geq x_1 \geq \dots \geq x_i\}$  for small enough  $\epsilon$ . Therefore,

$$\begin{aligned}
&\int_0^\epsilon dx_1 \int_0^{x_1} dx_2 \dots \int_0^{x_{i-1}} dx_i f^{(i)}(x_i) \\
&= \int f^{(i)}(x_i) \mathbb{I}(\vec{x} \in \mathbb{A}) d\vec{x} \\
&\leq C |\mathbb{A}| \\
&= \Theta(\epsilon^i).
\end{aligned}$$

□

As a corollary,

**Lemma C.19.** *If  $f : (0, 2) \rightarrow \mathbb{R}$  is smooth on  $(0, 2)$  and  $\lim_{x \rightarrow 0} f^{(i)}(x)$  exists and is finite for all  $i$ , then  $f$  can be extended to  $[0, 2)$  and be one-sided smooth at 0, and*

$$f(\epsilon) = f(0) + \epsilon f^{(1)}(0) + \dots + \frac{\epsilon^{i-1}}{(i-1)!} f^{(i-1)}(0) + O(\epsilon^i)$$

for any  $i$ .

**Lemma C.20.** *Let  $\phi = \tanh$ . For any fixed  $c$ ,  $W\phi(q, cq)$  is smooth (infinitely differentiable) on  $q \in (0, \infty)$ . As a function of  $Q := q^{-1}$ , it can be extended smoothly to the point  $Q = 0$ , so that*

$$\begin{aligned}
W\phi(q, cq) &= \lim_{q' \rightarrow \infty} W\phi(q', cq') + q^{-1} \lim_{q' \rightarrow \infty} \partial W\phi(q', cq') / \partial (q')^{-1} + \dots \\
&\quad + \frac{q^{-i+1}}{(i-1)!} \lim_{q' \rightarrow \infty} \partial^{i-1} W\phi(q', cq') / \partial (q')^{-i+1} + O(q^{-i})
\end{aligned}$$

for any  $i \geq 0$ . Furthermore, for  $c$  bounded away from 1, the constants hidden  $O$  can be taken independent of  $c$ .

*Proof. Smoothness on  $(0, \infty)$ .* By the third equation of [Lemma C.16](#), for  $Q \in (0, \infty) \iff q \in (0, \infty)$ ,

$$\begin{aligned}
&\frac{1}{\pi} \int_0^\infty s ds \left| \frac{\partial^n}{\partial Q^n} \left( Q e^{-s^2 Q/2} \right) \right| \int_0^\pi d\theta |\phi(s \sin \theta) \phi(s \sin(\theta - \tau))| \\
&\leq \int_0^\infty s ds \left| \frac{\partial^n}{\partial Q^n} \left( Q e^{-s^2 Q/2} \right) \right| < \infty,
\end{aligned}$$

so by Leibniz's integral rule and a simple induction, all derivatives of  $W\phi(q, cq)$  against  $Q$  exists for any  $Q \in (0, \infty)$ .

**Extension to  $Q = 0$ .** By [Lemma C.19](#), it suffices to show that the limit of  $\frac{\partial^k W\phi(q, cq)}{\partial Q^k}$  exists and is finite as  $Q \rightarrow 0$ , for all  $k$ . Let  $\tau = \arccos c$ . By the fourth equation of [Lemma C.16](#), we have explicitly

$$\begin{aligned}\frac{\partial^k W\phi(q, cq)}{\partial Q^k} &= \frac{1}{\pi} \int_0^\pi d\theta \int_0^\infty ds (-s^2/2)^k e^{-s^2 Q/2} \frac{\partial}{\partial s} \Sigma(s, \theta; \tau) \\ &= \frac{(-2)^{-k}}{\pi} \int_0^\pi d\theta \int_0^\infty ds s^{2k} e^{-s^2 Q/2} \frac{\partial}{\partial s} \Sigma(s, \theta; \tau)\end{aligned}$$

for any  $Q \in (0, \infty)$ . Note that for  $\phi = \tanh$ ,  $\dot{\phi} = \text{sech}^2$ ,

$$\frac{\partial}{\partial s} \Sigma(s, \theta; \tau) = \sin \theta \dot{\phi}(s \sin \theta) \phi(s \sin(\theta - \tau)) + \sin(\theta - \tau) \phi(s \sin \theta) \dot{\phi}(s \sin(\theta - \tau)).$$

We split the integral of  $\frac{\partial^k W\phi}{\partial Q^k}$  as follows:

$$\begin{aligned}\frac{\partial^k W\phi(q, cq)}{\partial Q^k} &= \frac{(-2)^{-k}}{\pi} \int_0^\pi d\theta \int_0^\infty ds s^{2k} e^{-s^2 Q/2} \sin \theta \dot{\phi}(s \sin \theta) \phi(s \sin(\theta - \tau)) \\ &\quad + \frac{(-2)^{-k}}{\pi} \int_0^\pi d\theta \int_0^\infty ds s^{2k} e^{-s^2 Q/2} \sin(\theta - \tau) \phi(s \sin \theta) \dot{\phi}(s \sin(\theta - \tau))\end{aligned}$$

We show that for each piece, the limit as  $Q \rightarrow 0$  exists and is finite, for any  $k$ . This will prove the smooth extendability of  $W\phi$  to  $Q = 0$ . We will do this for the first piece; the second is similar.

For  $Q > 0$ , the integrand is absolutely integrable, so we may switch the integrals.

$$\begin{aligned}&\int_0^\pi d\theta \int_0^\infty ds s^{2k} e^{-s^2 Q/2} \sin \theta \dot{\phi}(s \sin \theta) \phi(s \sin(\theta - \tau)) \\ &= \int_0^\infty ds s^{2k} e^{-s^2 Q/2} \int_0^\pi d\theta \sin \theta \dot{\phi}(s \sin \theta) \phi(s \sin(\theta - \tau))\end{aligned}$$

We now try to bound the inner integral by an exponentially decreasing term  $e^{-s\mu}$  for some  $\mu$ ; clearly, by monotone convergence on the outer integral as  $Q \rightarrow 0$ , this would show the limit of the integral exists and is finite.

Because  $\phi$  is odd and  $\dot{\phi}$  is even, the inner integrand is negative on  $\theta \in [0, \tau]$  and positive on  $\theta \in (\tau, \pi]$ . We will break up the inner integral as follows, for some fixed  $\epsilon > 0$  satisfying  $\tau - \epsilon > 0$  independent of  $s$  (recall  $\tau \in (0, \pi/2]$ ).

$$\begin{aligned}&\int_0^\pi d\theta \sin \theta \dot{\phi}(s \sin \theta) \phi(s \sin(\theta - \tau)) \\ &= \left( \int_0^\epsilon + \int_{\pi-\epsilon}^\pi \right) d\theta \sin \theta \dot{\phi}(s \sin \theta) \phi(s \sin(\theta - \tau)) + \int_\epsilon^{\pi-\epsilon} d\theta \sin \theta \dot{\phi}(s \sin \theta) \phi(s \sin(\theta - \tau))\end{aligned}$$

Now because  $\dot{\phi}(z) = \text{sech}^2(z) \leq 2e^{-z}$ , and  $\sin \theta \geq \sin \epsilon$  on  $\theta \in [\epsilon, \pi - \epsilon]$ ,

$$\begin{aligned}&\left| \int_\epsilon^{\pi-\epsilon} d\theta \sin \theta \dot{\phi}(s \sin \theta) \phi(s \sin(\theta - \tau)) \right| \\ &\leq 2 \int_\epsilon^{\pi-\epsilon} d\theta \exp(-s \sin \epsilon) \\ &= 2(\pi - 2\epsilon) \exp(-s \sin \epsilon).\end{aligned}$$

For the other part:

$$\begin{aligned}&\int_{\pi-\epsilon}^\pi d\theta \sin \theta \dot{\phi}(s \sin \theta) \phi(s \sin(\theta - \tau)) \\ &= \int_\epsilon^0 \sin(\pi - \theta) \dot{\phi}(s \sin \pi - \theta) \phi(s \sin(\pi - \theta - \tau)) d(\pi - \theta) \\ &= \int_0^\epsilon d\theta \sin \theta \dot{\phi}(s \sin \theta) \phi(s \sin \theta + \tau)\end{aligned}$$



so that

$$\begin{aligned} & \left( \int_0^\epsilon + \int_{\pi-\epsilon}^\pi \right) d\theta \sin \theta \dot{\phi}(s \sin \theta) \phi(s \sin(\theta - \tau)) \\ &= \int_0^\epsilon d\theta \sin \theta \dot{\phi}(s \sin \theta) [\phi(s \sin(\tau + \theta)) - \phi(s \sin(\tau - \theta))] \end{aligned}$$

But by intermediate value theorem,  $\phi(s \sin(\tau + \theta)) - \phi(s \sin(\tau - \theta)) = 2\theta \partial \phi(s \sin(\tau + \theta)) / \partial \theta|_{\theta=\psi} = 2\theta \dot{\phi}(s \sin(\tau + \psi)) s \cos(\tau + \psi)$  for some  $\psi \in [-\theta, \theta]$ . By the assumption on  $\epsilon$ ,  $\phi(s \sin(\tau + \theta)) - \phi(s \sin(\tau - \theta)) \leq 2\epsilon \dot{\phi}(s \sin(\tau - \epsilon)) s \cos(\tau - \epsilon)$ . Then

$$\begin{aligned} & \int_0^\epsilon d\theta \sin \theta \dot{\phi}(s \sin \theta) [\phi(s \sin(\tau + \theta)) - \phi(s \sin(\tau - \theta))] \\ & \leq \int_0^\epsilon d\theta \sin \theta \dot{\phi}(s \sin \theta) 2\epsilon \dot{\phi}(s \sin(\tau - \epsilon)) s \cos(\tau - \epsilon) \\ & \leq 2\epsilon \dot{\phi}(s \sin(\tau - \epsilon)) s \cos(\tau - \epsilon) O(1) \end{aligned}$$

Because  $\tau - \epsilon > 0$  by assumption on  $\epsilon$ , and because  $\dot{\phi}(z) = \exp(-\Theta_+(z))$ , this quantity is  $\exp(-\Theta_+(z))$ , as desired (here  $\Theta_+$  denotes a positive quantity).

Thus

$$\begin{aligned} & \int_0^\pi d\theta \sin \theta \dot{\phi}(s \sin \theta) \phi(s \sin(\theta - \tau)) \\ &= \left( \int_0^\epsilon + \int_{\pi-\epsilon}^\pi \right) d\theta \sin \theta \dot{\phi}(s \sin \theta) \phi(s \sin(\theta - \tau)) + \int_\epsilon^{\pi-\epsilon} d\theta \sin \theta \dot{\phi}(s \sin \theta) \phi(s \sin(\theta - \tau)) \\ &= \exp(-\Theta_+(s)) \end{aligned}$$

and similarly for the other piece of  $\frac{\partial^k W \phi}{\partial Q^k}$ , so that

$$\begin{aligned} & \int_0^\infty ds s^{2k} e^{-s^2 Q/2} \int_0^\pi d\theta \sin \theta \dot{\phi}(s \sin \theta) \phi(s \sin(\theta - \tau)) \\ &= \int_0^\infty ds s^{2k} e^{-s^2 \frac{Q}{2} - \Theta_+(z)} \\ & \rightarrow \int_0^\infty ds s^{2k} e^{-\Theta_+(z)} \end{aligned}$$

is finite as  $Q \rightarrow 0$ , by monotone convergence.

**Independence of constant hidden in  $O((q')^{-i})$ .** The constant hidden is a function of the  $\epsilon$  chosen above, which depend on  $\tau$ , but only to the extent that it must satisfy  $\tau - \epsilon > 0$ . As long as we are interested in a set  $\mathcal{C}$  of  $c$  that is bounded away from 1, the corresponding set of  $\tau$  is bounded away from 0, so  $\epsilon$  can be taken to be some number smaller than all of the corresponding  $\tau$ .

□

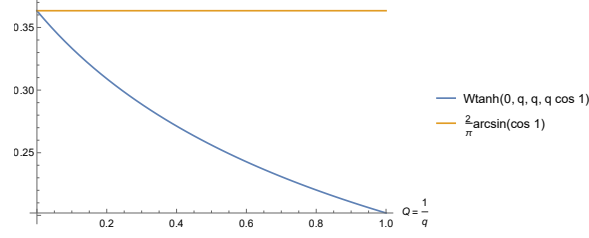
**Lemma C.21.** *Suppose  $\phi$  is tanh-like. Then for  $c \in [0, 1]$ ,*

$$W\phi(q, cq) \leq \frac{2}{\pi} \arcsin(c),$$

*and weakly increases to this upper bound as  $q \rightarrow \infty$ . Furthermore,*

- *If  $c = 0$  or 1, then equality holds regardless of  $q$ .*
- *If  $c \in (0, 1)$  is held constant,  $\frac{2}{\pi} \arcsin(c) - W\phi(q, cq) = \Theta(q^{-1})$ , where the hidden constants in  $\Theta$  depend on  $c$ . But the constants can be made independent of  $c$  if  $c \in [\epsilon, 1 - \epsilon]$  for some  $\epsilon > 0$ .*

*Proof.* The cases of  $c = 0$  or 1 are obvious by the definition of  $W$ . So from here on we assume  $c \in (0, 1)$ .



**Figure C.11:** We verify empirically that the subleading term in  $W \tanh(q, cq)$  is linear in  $q^{-1}$ , for constant  $c$ . Indeed, observe that the curve of  $W \tanh$  intersects the  $y$ -axis at an angle.

Let  $\tau := \arccos c$ . By the first equation of [Lemma C.16](#) and the assumption that  $\phi$  is tanh-like, it is immediate that  $W\phi(q, cq)$  is nondecreasing in  $q$ . By dominated convergence, using the second equation of [Lemma C.16](#), we get

$$\begin{aligned} \lim_{q \rightarrow \infty} W\phi(q, cq) &= \frac{1}{\pi} \int_0^\infty r \, dr e^{-r^2/2} (\pi - 2\tau) \\ &= \frac{\pi - 2\tau}{\pi} \\ &= \frac{2}{\pi} \arcsin c. \end{aligned}$$

Then the convergence rate is  $O(q^{-1})$  by [Lemma C.20](#) and Taylor's theorem. Thus to show the convergence rate is  $\Theta(q^{-1})$ , it suffices to show that  $\mathbf{D} := \frac{\partial W\phi(q, cq)}{\partial Q} < 0$ . But this is apparent from the first equation of [Lemma C.16](#): For  $\tau \in (0, \pi/2)$ ,

$$\begin{aligned} \mathbf{D} &= \frac{1}{\pi \sin \tau} \int_{w' \geq |w|} dw \, dw' \Upsilon(w, w'; \tau) \left( -\frac{1}{2\sqrt{2}} Q^{-3/2} \right) \\ &\quad \times [\dot{\phi}(\sqrt{q/2}(w + w')) \phi(\sqrt{q/2}(w' - w)) \\ &\quad + \phi(\sqrt{q/2}(w + w')) \dot{\phi}(\sqrt{q/2}(w' - w))] \\ &< 0 \end{aligned}$$

since  $\Upsilon$  is positive on the integration domain, and  $\dot{\phi}$  and  $\phi$  are both positive for positive arguments, by the assumption of  $\phi$  being tanh-like.

**Independence of the constants in  $\Theta(q^{-1})$  from  $c$  when  $c \in [\epsilon, 1 - \epsilon]$ .** By [Lemma C.20](#), the upper constant can be made independent from  $c$ . Since  $\mathbf{D}$  is monotonically decreasing in  $c$  (or monotonically increasing in  $\tau$ ) and  $|\mathbf{D}|$  is monotonically increasing in  $c$  (or monotonically decreasing in  $\tau$ ), we have

$$|\mathbf{D}| > |\mathbf{D}| \Big|_{c=\epsilon}, \text{ which can be taken to be the lower constant in } \Theta(q^{-1}). \quad \square$$

[Fig. C.11](#) verifies empirically that the subleading term in  $W \tanh(q, cq)$  is linear in  $q^{-1}$ , for constant  $c$ .

**Theorem B.4.** Suppose  $\phi$  is a tanh-like nonlinearity in an RRN. Assume  $\mathbf{e}^{(0)} < 1$ .

- If  $\sigma_w = 0$ , then  $\gamma^{(l)} = lW\phi(\sigma_b^2, \sigma_b^2) + \gamma^{(0)} = lV\phi(\sigma_b^2) + \gamma^{(0)}$  and  $\lambda^{(l)} = \sigma_b^2$ , so that  $\mathbf{e}^{(l)} \rightarrow 1$  and  $1 - \mathbf{e}^{(l)} = \Theta(l^{-1})$ . As a result,  $\mathbf{s}^{(l)} = \mathbf{p}^{(l)}(1 - \mathbf{e}^{(l)}) = \Theta(1)$ .
- If  $\sigma_w > 0$ , then  $\gamma^{(l)} = \check{\Theta}(l^{\frac{2}{\pi}})$ , and  $\mathbf{e}^{(l)} \rightarrow 0$  like  $\check{\Theta}(l^{\frac{2}{\pi}-1})$ . Thus  $\mathbf{s}^{(l)} = \Theta(\mathbf{p}^{(l)}) = \Theta(l)$ .

*Proof.* We have by [Lemma C.21](#),

$$\gamma = \frac{2}{\pi} \arcsin(\lambda/\mathbf{q}) - \Theta(\mathbf{q}^{-1}) + \underline{\gamma}.$$

Since  $\mathbf{q} = \sigma_w^2 \underline{\mathbf{p}} + \sigma_b^2$  by [Thm B.2](#), and  $\lambda = \sigma_w^2 \underline{\gamma} + \sigma_b^2$  by [Thm B.3](#),

$$\gamma = \frac{2}{\pi} \arcsin \left( \frac{\sigma_w^2 \underline{\gamma} + \sigma_b^2}{\sigma_w^2 \underline{\mathbf{p}} + \sigma_b^2} \right) - \Theta(\mathbf{q}^{-1}) + \underline{\gamma}.$$

We claim that  $\gamma^{(l)} \rightarrow \infty$  as  $l \rightarrow \infty$ . Otherwise, there is some  $C$  such that  $\gamma^{(l)} \leq C$  for all  $l$ . For large enough  $l$ ,  $\mathbf{p}^{(l)} \geq \omega l$  for any  $\omega < 1$  and  $\arcsin \left( \frac{C}{\sigma_w^2 \mathbf{p}^{(l-1)} + \sigma_b^2} \right) = \Theta(1/l)$  by linearization of  $\arcsin$ . Thus  $\gamma^{(l)} = \Theta(\log l)$ , but this contradicts our assumption that  $\gamma$  is bounded. This proves our claim.

Therefore, for large enough  $l$ ,

$$\frac{\sigma_w^2 \underline{\gamma} + \sigma_b^2}{\sigma_w^2 \underline{\mathbf{p}} + \sigma_b^2} = \underline{\gamma} / \underline{\mathbf{p}} + \Theta(l^{-1}).$$

[Fig. C.12](#) shows  $\frac{2}{\pi} \arcsin x$  vs  $x$ . One sees that 1 is an unstable fixed point; if  $e < 1 - \epsilon$ , then  $\frac{2}{\pi} \arcsin e < 1 - \epsilon - \delta$  for some  $\delta$ . Thus  $c$  drops monotonically until some threshold under which the linearization of  $\arcsin$ ,  $\arcsin x = x + \Theta(x^3)$ , is applicable. So for large enough  $l$ ,

$$\begin{aligned} \gamma - \underline{\gamma} &= \frac{2}{\pi} \arcsin(\underline{\gamma} / \underline{\mathbf{p}} + \Theta(l^{-1})) - \Theta(l^{-1}) \\ &= \frac{2}{\pi} \underline{\gamma} / \underline{\mathbf{p}} + O(l^{-1}) \end{aligned}$$

As  $\mathbf{p}^{(l)} \sim l$  by [Thm B.2](#), this difference equation has solution  $\gamma = \Omega(l^{\frac{2}{\pi} - \epsilon})$ ,  $O(l^{\frac{2}{\pi} + \epsilon})$  for any  $\epsilon$  by using the dynamics of [Lemma C.11](#) to upper and lower bound this difference equation.  $\square$

## C.5.2 Backward Dynamics

**Theorem B.6.** For  $\phi = \tanh$  in an RRN,

- If  $\sigma_w = 0$ ,  $\chi^{(m)} = \chi^{(l)}$  for all  $l, m$ .
- If  $\sigma_w > 0$ ,

$$\log(\chi^{(m)} / \chi^{(l)}) = \mathcal{A}(\sqrt{l} - \sqrt{m}) + \mathcal{B}(\log l - \log m) + O(1)$$

$$\text{where } \mathcal{A} = \frac{4}{3} \sqrt{\frac{2}{\pi}} \sigma_w \text{ and } \mathcal{B} = \frac{4}{3\pi} - \sigma_w^2 \frac{4}{9\pi}.$$

*Proof.* The  $\sigma_w = 0$  case is obvious. We will assume  $\sigma_w > 0$  from here on.

Let  $\mathbf{p}^{(l)} = b_0 l + b_1 l^{1/2} + b_2 \log l + O(1)$ . Then for  $D = \frac{2}{3} \sqrt{\frac{2}{\pi}}$ , we have (implicitly applying [Lemma C.4](#) and [Lemma C.8](#)),

$$\begin{aligned} \mathbf{q}^{-1/2} &= \sigma_w^{-1} b_0^{-1/2} l^{-1/2} (1 - b_1 b_0^{-1} 2^{-1} l^{-1/2} - b_2 b_0^{-1} 2^{-1} l^{-1} \log l + O(l^{-1})) \\ V\dot{\phi}(\mathbf{q}) &= D \mathbf{q}^{-1/2} + \Theta(\mathbf{q}^{-3/2}) \\ &= D \sigma_w^{-1} b_0^{-1/2} l^{-1/2} (1 - b_1 b_0^{-1} 2^{-1} l^{-1/2} - b_2 b_0^{-1} 2^{-1} l^{-1} \log l + O(l^{-1})) \\ \log(BV\dot{\phi}(\mathbf{q}) + 1) &= BD \sigma_w^{-1} b_0^{-1/2} l^{-1/2} \\ &\quad - (BD \sigma_w^{-1} b_0^{-3/2} b_1 2^{-1} + B^2 D^2 \sigma_w^{-2} b_0^{-1} 2^{-1}) l^{-1} + \Theta(l^{-3/2} \log l) \\ \sum_{r=1}^l \log(BV\dot{\phi}(\mathbf{q}^{(r)}) + 1) &= 2BD \sigma_w^{-1} b_0^{-1/2} l^{1/2} \\ &\quad - (BD \sigma_w^{-1} b_0^{-3/2} b_1 2^{-1} + B^2 D^2 \sigma_w^{-2} b_0^{-1} 2^{-1}) \log l + O(1) \end{aligned}$$

In our case, we have  $b_0 = 1, b_1 = -2C\sigma_w^{-1}, b_2 = C^2\sigma_w^{-2}, B = \sigma_w^2, C = \sqrt{\frac{2}{\pi}}$ , which gives

$$\sum_{r=1}^l \log(BV\dot{\phi}(\mathbf{q}^{(r)}) + 1) = \frac{4}{3}\sqrt{\frac{2}{\pi}}\sigma_w l^{1/2} + \left(\frac{4}{3\pi} - \sigma_w^2 \frac{4}{9\pi}\right) \log l + O(1).$$

so that

$$\chi^{(m)}/\chi^{(l)} = \exp \left[ \frac{4}{3}\sqrt{\frac{2}{\pi}}\sigma_w(\sqrt{l} - \sqrt{m}) + \left(\frac{4}{3\pi} - \sigma_w^2 \frac{4}{9\pi}\right)(\log l - \log m) + O(1) \right]$$

□

**Theorem B.7.** *Suppose  $\phi = \tanh$ . Then in an RRN*

- If  $\sigma_w = 0$ ,  $\chi_b^{(l)} = \chi^{(L)}V\dot{\phi}(\sigma_b^2)$  and  $\chi_w^{(l)} = \chi^{(L)}V\dot{\phi}(\sigma_b^2)((l-1)V\dot{\phi}(\sigma_b^2) + \mathbf{p}^{(0)})$ , where  $L$  is the last layer.
- If  $\sigma_w > 0$ ,

$$\log(\chi_b^{(m)}/\chi_b^{(l)}) = \mathcal{A}(\sqrt{l} - \sqrt{m}) + \mathcal{B}_b(\log l - \log m) + O(1)$$

$$\log(\chi_w^{(m)}/\chi_w^{(l)}) = \mathcal{A}(\sqrt{l} - \sqrt{m}) + \mathcal{B}_w(\log l - \log m) + O(1)$$

where  $\mathcal{A} = \frac{4}{3}\sqrt{\frac{2}{\pi}}\sigma_w$  (same as  $\mathcal{A}$  in [Thm B.6](#)) and  $\mathcal{B}_b = \mathcal{B} + \frac{1}{2}, \mathcal{B}_w = \mathcal{B} - \frac{1}{2}$ , with  $\mathcal{B} = \frac{4}{3\pi} - \sigma_w^2 \frac{4}{9\pi}$  (same as  $\mathcal{B}$  in [Thm B.6](#)).

*Proof.* The  $\sigma_w = 0$  case is obvious. We will assume  $\sigma_w > 0$  from here on.

As in the proof of [Thm B.6](#),

$$V\dot{\phi}(\mathbf{q}) = D\sigma_w^{-1}b_0^{-1/2}l^{-1/2} + \Theta(l^{-1})$$

where  $D = \frac{2}{3}\sqrt{\frac{2}{\pi}}$ . Thus by [Thm B.5](#),

$$\log(\chi^{(m)}/\chi^{(l)}) = \frac{4}{3}\sqrt{\frac{2}{\pi}}\sigma_w(\sqrt{l} - \sqrt{m}) + \left(\frac{4}{3\pi} - \sigma_w^2 \frac{4}{9\pi}\right)(\log l - \log m) + O(1)$$

$$\log(\chi_b^{(m)}/\chi_b^{(l)}) = \frac{4}{3}\sqrt{\frac{2}{\pi}}\sigma_w(\sqrt{l} - \sqrt{m}) + \left(\frac{4}{3\pi} - \frac{1}{2} - \sigma_w^2 \frac{4}{9\pi}\right)(\log l - \log m) + O(1)$$

Similarly, since  $\mathbf{p} = l + \Theta(\sqrt{l})$  by [Thm B.2](#), we have

$$\log(\chi_w^{(m)}/\chi_w^{(l)}) = \frac{4}{3}\sqrt{\frac{2}{\pi}}\sigma_w(\sqrt{l} - \sqrt{m}) + \left(\frac{4}{3\pi} + \frac{1}{2} - \sigma_w^2 \frac{4}{9\pi}\right)(\log l - \log m) + O(1)$$

□

## C.6 Tanh: Full Residual Network

### C.6.1 Forward Dynamics

**Theorem B.9.** *Suppose  $\phi$  is tanh-like. Assume the FRN architecture.*

- If  $\sigma_w = 0$ , then  $\mathbf{p}^{(l)} = (\sigma_v^2 V\dot{\phi}(\sigma_b^2) + \sigma_a^2)l + \mathbf{p}^{(0)}$ , and  $\mathbf{q}^{(l)} = \sigma_b^2$ .
- If  $\sigma_w > 0$ , then  $\mathbf{p}^{(l)} = b_0 l + b_1 l^{1/2} + b_2 \log l + O(1)$ , where

$$\begin{aligned} b_0 &= \sigma_v^2 + \sigma_a^2 \\ b_1 &= \frac{-2C\sigma_v^2\sigma_w^{-1}}{\sqrt{\sigma_v^2 + \sigma_a^2}} \\ b_2 &= \frac{-C^2\sigma_v^4\sigma_w^{-2}}{(\sigma_v^2 + \sigma_a^2)^2} \end{aligned}$$

and  $C = \sqrt{\frac{2}{\pi}}$ . Additionally,  $\mathbf{q}^{(l)} = \sigma_w^2 b_0 l + \sigma_w^2 b_1 l^{1/2} + \sigma_w^2 b_2 \log l + O(1)$ .

*Proof.* The  $\sigma_w = 0$  case is obvious. We will assume  $\sigma_w > 0$  from here on.

As in **Thm B.2**,  $\mathbf{p}$  will have expansion  $\mathbf{p} = b_0 l + b_1 l^{1/2} + b_2 \log l + O(1)$ . Then, for  $C = \sqrt{\frac{2}{\pi}}$ ,

$$\begin{aligned} \mathbf{q}^{-1/2} &= \sigma_w^{-1} b_0^{-1/2} l^{-1/2} (1 - b_1 b_0^{-1} 2^{-1} l^{-1/2} - b_2 b_0^{-1} 2^{-1} l^{-1} \log l + O(l^{-1})) \\ \sum_{r=1}^l V\phi(\mathbf{q}^{(r)}) &= \sum_{r=1}^l 1 - C(\mathbf{q}^{(r)})^{-1/2} + \Theta((\mathbf{q}^{(r)})^{-3/2}) \\ &= l - 2C\sigma_w^{-1} b_0^{-1/2} l^{1/2} + C\sigma_w^{-1} b_1 b_0^{-3/2} 2^{-1} \log l + O(1) \\ \mathbf{p}^{(l)} &= \sigma_v^2 \sum_{r=1}^l + \sigma_a^2 l \\ &= (\sigma_v^2 + \sigma_a^2) l - 2C\sigma_v^2 \sigma_w^{-1} b_0^{-1/2} l^{1/2} + C\sigma_v^2 \sigma_w^{-1} b_1 b_0^{-3/2} 2^{-1} \log l + O(1) \end{aligned}$$

which yields

$$\begin{aligned} b_0 &= \sigma_v^2 + \sigma_a^2 \\ b_1 &= -2C\sigma_v^2 \sigma_w^{-1} b_0^{-1/2} = \frac{-2C\sigma_v^2 \sigma_w^{-1}}{\sqrt{\sigma_v^2 + \sigma_a^2}} \\ b_2 &= \frac{-C^2 \sigma_v^4 \sigma_w^{-2}}{(\sigma_v^2 + \sigma_a^2)^2} \end{aligned}$$

□

**Lemma C.22.** *Suppose  $\phi$  is tanh-like. Then*

$$\gamma \leq \sigma_v^2 \frac{2}{\pi} \arcsin(\lambda/\mathbf{q}) + \sigma_a^2 + \underline{\gamma},$$

and

$$\sigma_v^2 \frac{2}{\pi} \arcsin(\lambda/\mathbf{q}) + \sigma_a^2 + \underline{\gamma} - \gamma = \Theta(\mathbf{q}^{-1}).$$

*Proof.* Similar to the proof of **Lemma C.21**. □

**Lemma C.23.** *Let  $u^* \in [0, 1)$ . Let  $f_t : [0, 1) \rightarrow [0, 1]$  be a continuous function for each  $t \in \mathbb{N}$ , to each of which we associate two numbers  $0 \leq a_t \leq u^* \leq b_t \leq 1$ . Suppose for each  $t$ ,  $f_t(u) > u$  for all  $u \in [0, a_t)$  and  $f_t(u) < u$  for all  $u \in (b_t, 1)$ . Assume that for each  $u$ ,  $f_t(u) - u \rightarrow 0$  as  $t \rightarrow \infty$  uniformly over  $u$ . If  $a_t \nearrow u^*$  and  $b_t \searrow u^*$ , then for any  $u_0 \in [0, 1)$ , the dynamics  $u_t = f_t(u_{t-1})$  has a limiting point. Furthermore, either  $u_t \rightarrow u^*$  or  $u_t$  eventually converges monotonically (decreasing or increasing) to a limit point.*

*Proof.* Fix a  $u_0 \in [0, 1)$ . If  $u_t \rightarrow u^*$  then we are done. Otherwise, suppose there is a neighborhood  $[u^* - \epsilon, u^* + \epsilon]$  such that for an infinite sequence  $t_1, t_2, \dots, u_{t_i} \notin [u^* - \epsilon, u^* + \epsilon]$ . WLOG assume  $u_{t_i} < u^* - \epsilon$  for all  $i$  and  $(t_i)_i$  is the sequence of all  $t$ s that satisfy this inequality.

If  $(t_i)_i$  contains  $\{s : s \geq N\}$  for some  $N$ , then for some  $M > N$ , for every  $t > M$ ,  $a_t > u^* - \epsilon > u_t$ . By assumption,  $u_t$  is monotonic for all  $t > M$  but is bounded above. Thus  $u_t$  has a fixed point  $\hat{u} \leq u^* - \epsilon$  as desired.

Now assume there are infinite  $i$ s such that  $t_i - 1 \neq t_{i-1}$  (i.e.  $t_i - 1$  is not part of the sequence  $(t_i)_i$ ). We will show that this case is contradictory. Take  $T$  large enough such that  $a_t > u^* - \epsilon/2$  and  $|f_t(u) - u| < \epsilon/4$  for all  $u$  and for all  $t \geq T$  ( $T$  exists by premise). Let  $j$  be the smallest index such that  $t_j > T$  and  $t_j - 1 \neq t_{j-1}$ . By the definition of  $j$ ,  $u_{t_j-1} \geq u^* - \epsilon$ . If  $u_{t_j-1} \geq u^* - \epsilon/2$ , then by definition of  $T$ ,  $u^* - \epsilon > u_{t_j} = f_{t_j}(u_{t_j-1}) > u_{t_j-1} - \epsilon/4 > u^* - 3\epsilon/4 > u^* - \epsilon$ , a contradiction. If  $u^* - \epsilon \leq u_{t_j-1} < u^* - \epsilon/2$ , then by the definition of  $T$ ,  $u_{t_j-1} \leq a_{t_j-1}$  so that  $u_{t_j} = f_{t_j}(u_{t_j-1}) > u_{t_j-1} \geq u^* - \epsilon$ , a contradiction.

The “furthermore” claim is clear from our proof above. □

**Theorem B.11.** *Assume  $\phi = \tanh$  in an FRN. Suppose  $\mathbf{e}^{(0)} < 1$ .*

- If  $\sigma_w = 0$ , then  $\lambda^{(l)} = \sigma_b^2$  and  $\gamma^{(l)} = l(\sigma_v^2 W\phi(\sigma_b^2, \sigma_b^2) + \sigma_a^2) + \gamma^{(0)} = l(\sigma_v^2 V\phi(\sigma_b^2) + \sigma_a^2) + \gamma^{(0)}$ . Thus  $\mathbf{e}^{(l)} \rightarrow 1$  and  $1 - \mathbf{e}^{(l)} = \Theta(l^{-1})$ . As a result,  $\mathbf{s}^{(l)} = \mathbf{p}^{(l)}(1 - \mathbf{e}^{(l)}) = \Theta(1)$ .
- If  $\sigma_w > 0$ , then  $\mathbf{e}^{(l)}$  converges to the unique fixed point  $\mathbf{e}^* \neq 1$  determined by the equation

$$\mathbf{e}^* = \frac{1}{\sigma_v^2 + \sigma_a^2} [\sigma_v^2 \frac{2}{\pi} \arcsin(\mathbf{e}^*) + \sigma_a^2].$$

Furthermore,  $\mathbf{e}^{(l)}$  converges to  $\mathbf{e}^*$  polynomially:  $|\mathbf{e}^{(l)} - \mathbf{e}^*|$  is  $\check{\Theta}(l^{-\delta^*})$ , where

$$\delta^* := 1 - \frac{2}{\pi} \frac{1}{\sqrt{1 - (\mathbf{e}^*)^2}} \frac{\sigma_v^2}{\sigma_v^2 + \sigma_a^2} \in \left[ \frac{2}{\pi} - 1, \frac{1}{2} \right)$$

Since  $\mathbf{e}^* < 1$ ,  $\mathbf{s}^{(l)} = \Theta(\mathbf{p}^{(l)}) = \Theta(l)$ .

*Proof.* The  $\sigma_w = 0$  case is obvious. We will assume  $\sigma_w > 0$  from here on.

If  $\sigma_a = 0$ , then  $\mathbf{e}^*$  as defined above is 0, and  $\mathbf{e} = \frac{\gamma}{\mathbf{p}}$  decreases as  $\Theta(l^{\frac{2}{\pi}-1})$  to 0, by the same reason as before.

So from now on suppose  $\sigma_a > 0$ . We apply [Lemma C.23](#) first to show that  $\mathbf{e}$  converges. We have

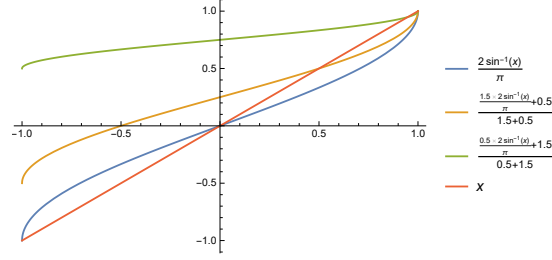
$$\begin{aligned} \sigma_v^2 W\phi(\mathbf{q}, c\mathbf{q}) + \sigma_a^2 &= \mathbf{e}\mathbf{p} - \underline{\mathbf{e}}\underline{\mathbf{p}} \\ &= \mathbf{e}\mathbf{p} - \underline{\mathbf{e}}\mathbf{p} + \underline{\mathbf{e}}\mathbf{p} - \underline{\mathbf{e}}\underline{\mathbf{p}} \\ &= (\mathbf{e} - \underline{\mathbf{e}})\mathbf{p} + \underline{\mathbf{e}}(\mathbf{p} - \underline{\mathbf{p}}) \\ &= (\mathbf{p} - \underline{\mathbf{p}})[(\mathbf{e} - \underline{\mathbf{e}}) \frac{\mathbf{p}}{\mathbf{p} - \underline{\mathbf{p}}} + \underline{\mathbf{e}}] \\ \frac{\sigma_v^2 W\phi(\mathbf{q}, c\mathbf{q}) + \sigma_a^2}{\sigma_v^2 V\phi(\mathbf{q}) + \sigma_a^2} &= (\mathbf{e} - \underline{\mathbf{e}}) \frac{\mathbf{p}}{\mathbf{p} - \underline{\mathbf{p}}} + \underline{\mathbf{e}} \\ \frac{\mathbf{p} - \underline{\mathbf{p}}}{\mathbf{p}} \left[ \frac{\sigma_v^2 W\phi + \sigma_a^2}{\sigma_v^2 V\phi + \sigma_a^2} - \underline{\mathbf{e}} \right] &= \mathbf{e} - \underline{\mathbf{e}} \end{aligned}$$

If we define  $f_l(u) := \frac{\mathbf{p}^{(l)} - \mathbf{p}^{(l-1)}}{\mathbf{p}^{(l)}} \left[ \frac{\sigma_v^2 W\phi(\mathbf{q}^{(l)}, c^{(l)}\mathbf{q}^{(l)}) + \sigma_a^2}{\sigma_v^2 V\phi(\mathbf{q}^{(l)}) + \sigma_a^2} - u \right] + u$  (the LHS of the above), then  $f_l(u) - u = O(l^{-1})$  uniformly for all  $u$  because  $\mathbf{p}^{(l)} = \Theta(l)$ ,  $\mathbf{p}^{(l)} - \mathbf{p}^{(l-1)} = \Theta(1)$ , and the part in the bracket is  $O(1)$ , with constants all (able to be taken) independent of  $u$ . We divide  $[0, 1)$  into the following intervals  $I_1 = [1, 1/2)$ ,  $I_2 = [1/2, 3/4)$ ,  $I_3 = [3/4, 7/8)$ ,  $\dots$ . For each  $I_k$ , it is clear that the trajectories of  $\mathbf{e}^{(l)} = f_l(\mathbf{e}^{(l-1)})$  with  $\mathbf{e}^{(0)} \in I_k$  will fall into some interval  $J_k$  bounded away from 1 for all  $l \geq L$ , for large enough  $L$  (dependent on  $k$ ). Then we can apply [Lemmas C.1](#), [C.5](#) and [C.21](#) to get  $f_l(u) = \frac{\mathbf{p}^{(l)} - \mathbf{p}^{(l-1)}}{\mathbf{p}^{(l)}} \left[ \frac{\sigma_v^2 \frac{2}{\pi} \arcsin(u) + \sigma_a^2}{\sigma_v^2 + \sigma_a^2} - u + o(1) \right] + u$  where the constants in  $o(1)$  is uniform for all  $\mathbf{e}^{(0)} \in I_k$ . For  $u < \mathbf{e}^*$  (as defined in the theorem statement),  $\frac{\sigma_v^2 \frac{2}{\pi} \arcsin(u) + \sigma_a^2}{\sigma_v^2 + \sigma_a^2} > u$  and for  $u > \mathbf{e}^*$ ,  $\frac{\sigma_v^2 \frac{2}{\pi} \arcsin(u) + \sigma_a^2}{\sigma_v^2 + \sigma_a^2} < u$  (see [Fig. C.12](#)). Thus as  $l \rightarrow \infty$ , the  $o(1)$  term gets smaller and smaller, and this monotonicity holds for  $f_l(u) - u = \left[ \frac{\sigma_v^2 \frac{2}{\pi} \arcsin(u) + \sigma_a^2}{\sigma_v^2 + \sigma_a^2} - u + o(1) \right] > 0$  (resp.  $< 0$ ) on larger and larger intervals  $[0, a_l] \cap J_k$  (resp.  $[b_l, 1) \cap J_k$ ). This proves all the preconditions for [Lemma C.23](#), which yields that  $I_k$  converges to a limit point. As this argument is independent of  $k$ , we have that for all  $\mathbf{e}^{(0)} \in [0, 1)$ ,  $\mathbf{e}^{(l)}$  converges.

Now we solve for the limit point.

Suppose  $\mathbf{e}$  has limit point  $\mathbf{e}^\dagger$  (possibly different from  $\mathbf{e}^*$  described in the theorem); if we express  $\gamma^{(l)} = (\mathbf{e}^\dagger + \epsilon^{(l)})\mathbf{p}^{(l)}$ , then

$$\begin{aligned} \sigma_v^2 W\phi(\mathbf{q}, c\mathbf{q}) + \sigma_a^2 &= \gamma - \underline{\gamma} \\ &= (\mathbf{e}^\dagger + \epsilon)\mathbf{p} - (\mathbf{e}^\dagger + \underline{\epsilon})\underline{\mathbf{p}} \\ &= \mathbf{e}^\dagger(\mathbf{p} - \underline{\mathbf{p}}) + \epsilon\mathbf{p} - \underline{\epsilon}\underline{\mathbf{p}} \\ \frac{\sigma_v^2 W\phi(\mathbf{q}, c\mathbf{q}) + \sigma_a^2}{\sigma_v^2 V\phi(\mathbf{q}) + \sigma_a^2} &= \mathbf{e}^\dagger + \epsilon + (\epsilon - \underline{\epsilon}) \frac{\mathbf{p}}{\mathbf{p} - \underline{\mathbf{p}}} \end{aligned}$$



**Figure C.12:** Graph of  $y(\mathbf{e}) = \frac{1}{\sigma_v^2 + \sigma_a^2} [\sigma_v^2 \frac{2}{\pi} \arcsin(\mathbf{e}) + \sigma_a^2]$  for various  $\sigma_v$  and  $\sigma_a$ .

As  $l \rightarrow \infty$ ,  $c \sim \mathbf{e} \rightarrow \mathbf{e}^\dagger$ , and  $W\phi(\mathbf{q}, \mathbf{e}^\dagger \mathbf{q}) \rightarrow \frac{2}{\pi} \arcsin(\mathbf{e}^\dagger)$ , and  $V\phi(\mathbf{q}) \rightarrow 1$ . Additionally,  $\underline{\mathbf{p}}/(\underline{\mathbf{p}} - \underline{\mathbf{p}}) = \Theta(l)$  and  $\epsilon = o(1)$  so that  $\epsilon - \underline{\epsilon} = o(l^{-1})$ . Then we have, taking limits  $l \rightarrow \infty$ ,

$$\frac{\sigma_v^2 \frac{2}{\pi} \arcsin(\mathbf{e}^\dagger) + \sigma_a^2}{\sigma_v^2 + \sigma_a^2} = \mathbf{e}^\dagger.$$

Since  $f_l$  (as defined above) repels points away from 1, the only solution for  $\mathbf{e}^\dagger$  when  $\mathbf{e}^{(0)} < 1$  is  $\mathbf{e}^\dagger = \mathbf{e}^*$  as specified in the theorem statement.

We defer the proof of the convergence rate to  $\mathbf{e}^*$  to [Thm C.25](#).

□

**Lemma C.24.** *Let  $\mathbf{e}^*$  be the stable fixed point determined by  $\sigma_a$  and  $\sigma_v$ . Then as long as  $\sigma_v > 0$ ,*

$$\frac{2}{\pi} \frac{1}{\sqrt{1 - (\mathbf{e}^*)^2}} \frac{\sigma_v^2}{\sigma_v^2 + \sigma_a^2} \in \left(\frac{1}{2}, \frac{2}{\pi}\right]$$

*Proof.* Write  $\rho := \frac{\sigma_a^2}{\sigma_v^2}$ . By definition of  $\mathbf{e}^*$ , we get

$$\begin{aligned} \mathbf{e}^* &= (1 - \rho) \frac{2}{\pi} \arcsin \mathbf{e}^* + \rho \\ \rho &= \frac{\mathbf{e}^* - \frac{2}{\pi} \arcsin \mathbf{e}^*}{1 - \mathbf{e}^*} \end{aligned}$$

Substituting  $\rho$  into the expression in question, it follows that we want to show

$$\frac{2}{\pi} (1 - \mathbf{e}^{*2})^{-1/2} (1 + \rho)^{-1} = \frac{2}{\pi} (1 - \mathbf{e}^{*2})^{-1/2} \left( \frac{1 - \frac{2}{\pi} \arcsin \mathbf{e}^*}{1 - \mathbf{e}^*} \right)^{-1} \in \left(\frac{1}{2}, \frac{2}{\pi}\right]$$

for  $\mathbf{e}^* \in [0, 1)$  (the endpoint at 1 is not included since  $\sigma_v > 0$ ). But this is

$$\frac{2}{\pi} (1 - \mathbf{e}^*)^{1/2} (1 + \mathbf{e}^*)^{-1/2} (1 - \frac{2}{\pi} \arcsin \mathbf{e}^*)^{-1}.$$

Set  $g(\mathbf{e}^*)$  to be this expression. We could proceed by finding critical points, but a simple plot [Fig. C.13](#) shows that  $g$  is decreasing on  $[0, 1)$ , with extremal values at the end points:

$$g(\mathbf{e}^*) \in \left[ \lim_{\mathbf{e}^* \rightarrow 1} g(\mathbf{e}^*), g(0) \right), \quad \text{for } \mathbf{e}^* \in [0, 1).$$

Obviously  $g(0) = \frac{2}{\pi}$ . For the limit, we note that  $\arcsin \mathbf{e}^*$  has an asymptotic expansion  $\frac{\pi}{2} - \sqrt{2}(1 - \mathbf{e}^*)^{1/2} + \Theta((1 - \mathbf{e}^*)^{3/2})$  at 1, so that  $(1 - \mathbf{e}^*)^{1/2} (1 - \frac{2}{\pi} \arcsin \mathbf{e}^*)^{-1} \rightarrow \frac{\pi}{2\sqrt{2}}$ , and  $g(\mathbf{e}^*) \rightarrow \frac{1}{2}$  as  $\mathbf{e}^* \rightarrow 1$ .

□

**Theorem C.25.** *If  $\mathbf{e}^{(0)} < 1$ , then  $|\mathbf{e}^{(l)} - \mathbf{e}^*|$  is  $\Omega(l^{-\delta^* - \epsilon})$  and  $O(l^{-\delta^* + \epsilon})$  for any  $\epsilon > 0$ , where*

$$\delta^* := 1 - \frac{2}{\pi} \frac{1}{\sqrt{1 - (\mathbf{e}^*)^2}} \frac{\sigma_v^2}{\sigma_v^2 + \sigma_a^2} \in \left[1 - \frac{2}{\pi}, \frac{1}{2}\right),$$

where the bounds on the right follow from [Lemma C.24](#).

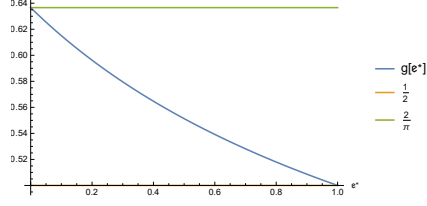


Figure C.13: Plot of  $g(\mathbf{e}^*)$  in the proof of Lemma C.24

*Proof.* Define  $\omega(q, c) = \frac{2}{\pi} \arcsin(c) - W \tanh(q, cq)$ . By Lemma C.21, for large enough  $l$ ,  $c$  is close to  $\mathbf{e}^*$  bounded away from 0 or 1, so that  $\omega(\mathbf{q}, c) = \Theta(\mathbf{q}^{-1})$  with the constant hidden in  $\Theta$  independent of  $c$ . Additionally, by Lemma C.5,  $1 - V \tanh(\mathbf{q}) = \Theta(\mathbf{q}^{-1/2})$ . Therefore,

$$\begin{aligned}
(\mathbf{e}^* + \epsilon)\mathbf{p} &= \sigma_v^2 \left( \frac{2}{\pi} \arcsin(\mathbf{e}^* + \epsilon) - \omega(\mathbf{q}, c) \right) + \sigma_a^2 + \underline{\gamma} \\
&= \sigma_v^2 \frac{2}{\pi} \left[ \arcsin(\mathbf{e}^*) + \frac{\epsilon}{\sqrt{1 - (\mathbf{e}^*)^2}} + \Theta(\epsilon^2) \right] - \Theta(l^{-1}) + \sigma_a^2 + \underline{\gamma} \\
&= \mathbf{e}^* (\sigma_v^2 + \sigma_a^2) + (\mathbf{e}^* + \epsilon)\underline{\mathbf{p}} + \sigma_v^2 \frac{2}{\pi} \frac{\epsilon}{\sqrt{1 - (\mathbf{e}^*)^2}} + \Theta(\epsilon^2) - \Theta(l^{-1}) \\
\mathbf{e}^* (\mathbf{p} - \underline{\mathbf{p}} - \sigma_v^2 - \sigma_a^2) &= \epsilon \underline{\mathbf{p}} - \epsilon \mathbf{p} + \sigma_v^2 \frac{2}{\pi} \frac{\epsilon}{\sqrt{1 - (\mathbf{e}^*)^2}} + \Theta(\epsilon^2) - \Theta(l^{-1}) \\
\mathbf{e}^* \sigma_v^2 (V\phi(\mathbf{q}) - 1) &= \epsilon \underline{\mathbf{p}} - \epsilon \mathbf{p} + \sigma_v^2 \frac{2}{\pi} \frac{\epsilon}{\sqrt{1 - (\mathbf{e}^*)^2}} + \Theta(\epsilon^2) - \Theta(l^{-1}) \\
\epsilon &= \frac{1}{\mathbf{p}} (\mathbf{e}^* \sigma_v^2 (1 - V\phi(\mathbf{q})) + \Theta(\epsilon^2) - \Theta(l^{-1}) + \epsilon (\underline{\mathbf{p}} + \sigma_v^2 \frac{2}{\pi} \frac{1}{\sqrt{1 - (\mathbf{e}^*)^2}})) \\
&= \Theta(l^{-3/2}) + \epsilon (1 - \delta^{(l)})/l
\end{aligned}$$

where

$$\begin{aligned}
\delta^{(l)} &= \frac{l}{\mathbf{p}} (\sigma_v^2 V\phi(\mathbf{q}) + \sigma_a^2 - \sigma_v^2 \frac{2}{\pi} \frac{1}{\sqrt{1 - (\mathbf{e}^*)^2}}) + \Theta(\epsilon/l) \\
&= (1 + \Theta(l^{-1/2})) (\sigma_v^2 (1 - \Theta(l^{-1/2})) + \sigma_a^2 - \sigma_v^2 \frac{2}{\pi} \frac{1}{\sqrt{1 - (\mathbf{e}^*)^2}}) / (\sigma_v^2 + \sigma_a^2) + \Theta(\epsilon/l) \\
&= \delta^* + O(l^{-1/2}),
\end{aligned}$$

where  $\delta^* := 1 - \frac{2}{\pi} \frac{1}{\sqrt{1 - (\mathbf{e}^*)^2}} \frac{\sigma_v^2}{\sigma_v^2 + \sigma_a^2}$ , which is positive by Lemma C.24. By taking the  $\delta$  of Lemma C.11 to be  $\delta^* + \epsilon$  or  $\delta^* - \epsilon$  respectively for lower and upper bounding the dynamics of  $\epsilon^{(l)}$ , the solution  $\epsilon^{(l)}$  is  $\Omega(l^{-\delta^* - \epsilon})$  and  $O(l^{-\delta^* + \epsilon})$  for any  $\epsilon > 0$  since  $\frac{1}{2} > \delta^*$ .  $\square$

## C.6.2 Backward Dynamics

**Theorem B.13.** Assume  $\phi = \tanh$  in an FRN.

- If  $\sigma_w = 0$ ,  $\chi^{(m)} = \chi^{(l)}$  for all  $l, m$ .
- If  $\sigma_w > 0$ , then for  $l \geq m \geq 0$ ,

$$\log(\chi^{(m)}/\chi^{(l)}) = \mathcal{A}(\sqrt{l} - \sqrt{m}) + \mathcal{B}(\log l - \log m) + O(1)$$



where

$$\mathcal{A} = \frac{4}{3} \sqrt{\frac{2}{\pi}} \frac{\sigma_v^2 \sigma_w}{\sqrt{\sigma_v^2 + \sigma_a^2}}$$

$$\mathcal{B} = \frac{4}{9\pi} \frac{\sigma_v^4}{\sigma_v^2 + \sigma_a^2} \left( \frac{3}{\sigma_v^2 + \sigma_a^2} - \sigma_w^2 \right)$$

*Proof.* The  $\sigma_w = 0$  case is obvious. We will assume  $\sigma_w > 0$  from here on.

As in the proof of [Thm B.6](#),

$$\begin{aligned} \log(\chi^{(m)}/\chi^{(l)}) &= 2BD\sigma_w^{-1}b_0^{-1/2}(\sqrt{l} - \sqrt{m}) \\ &\quad - (BD\sigma_w^{-1}b_0^{-3/2}b_12^{-1} + B^2D^2\sigma_w^{-2}b_0^{-1}2^{-1})(\log l - \log m) + O(1) \end{aligned}$$

where  $B = \sigma_v^2\sigma_w^2$ ,  $D = \frac{2}{3}\sqrt{\frac{2}{\pi}}$ ,

$$\begin{aligned} b_0 &= \sigma_v^2 + \sigma_a^2 \\ b_1 &= \frac{-2C\sigma_v^2\sigma_w^{-1}}{\sqrt{\sigma_v^2 + \sigma_a^2}} \\ b_2 &= \frac{-C^2\sigma_v^4\sigma_w^{-2}}{(\sigma_v^2 + \sigma_a^2)^2}. \end{aligned}$$

with  $C = \sqrt{\frac{2}{\pi}}$ . This simplifies to the desired form. □

**Theorem B.14.** Suppose  $\phi = \tanh$  in an FRN.

- If  $\sigma_w = 0$ , then

$$\begin{aligned} \chi_b^{(l)} &= \sigma_v^2 \chi^{(L)} \mathbf{V} \dot{\phi}(\sigma_b^2) \\ \chi_w^{(l)} &= \sigma_v^2 \chi^{(L)} \mathbf{V} \dot{\phi}(\sigma_b^2) ((\sigma_v^2 \mathbf{V} \phi(\sigma_b^2) + \sigma_a^2)(l-1) + \mathbf{p}^{(0)}) \\ \chi_v^{(l)} &= \chi^{(L)} \mathbf{V} \phi(\sigma_b^2) \\ \chi_a^{(l)} &= \chi^{(L)}. \end{aligned}$$

- If  $\sigma_w > 0$ , then for  $l \geq m \geq 0$ ,

$$\begin{aligned} \log(\chi_b^{(m)}/\chi_b^{(l)}) &= \mathcal{A}(\sqrt{l} - \sqrt{m}) + \mathcal{B}_b(\log l - \log m) + O(1) \\ \log(\chi_w^{(m)}/\chi_w^{(l)}) &= \mathcal{A}(\sqrt{l} - \sqrt{m}) + \mathcal{B}_w(\log l - \log m) + O(1) \\ \log(\chi_a^{(m)}/\chi_a^{(l)}) &= \mathcal{A}(\sqrt{l} - \sqrt{m}) + \mathcal{B}(\log l - \log m) + O(1) \\ \log(\chi_v^{(m)}/\chi_v^{(l)}) &= \mathcal{A}(\sqrt{l} - \sqrt{m}) + \mathcal{B}(\log l - \log m) + O(1) \end{aligned}$$

where  $\mathcal{A} = \frac{4}{3} \sqrt{\frac{2}{\pi}} \frac{\sigma_v^2 \sigma_w}{\sqrt{\sigma_v^2 + \sigma_a^2}}$  and  $\mathcal{B} = \frac{4}{9\pi} \frac{\sigma_v^4}{\sigma_v^2 + \sigma_a^2} \left( \frac{3}{\sigma_v^2 + \sigma_a^2} - \sigma_w^2 \right)$  are as in [Thm B.13](#) and  $\mathcal{B}_b = \mathcal{B} + \frac{1}{2}$  and  $\mathcal{B}_w = \mathcal{B} - \frac{1}{2}$ .

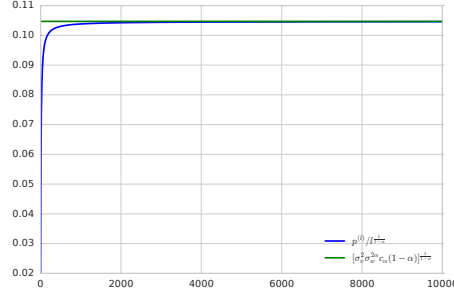
*Proof.* Similar to [Thm B.7](#). □

## C.7 $\alpha$ -ReLU: Full Residual Network

The following can be checked readily

**Lemma B.15.** If  $\alpha > -\frac{1}{2}$ , then  $\mathbf{V}\psi_\alpha(q) = c_\alpha q^\alpha$ , where  $c_\alpha = \frac{1}{\sqrt{\pi}} 2^{\alpha-1} \Gamma(\alpha + \frac{1}{2})$ .

Since  $\dot{\psi}_\alpha = \alpha \psi_{\alpha-1}$ , we have as a corollary,



**Figure C.14:** Verification of leading term of **Thm C.28** for  $\alpha = 0.55$ .

**Lemma C.26.** *If  $\alpha > \frac{1}{2}$ , then  $V\dot{\psi}_\alpha(q) = \alpha^2 c_{\alpha-1} q^{\alpha-1}$ .*

As a special case, when  $\alpha = 1$ ,  $c_\alpha = \frac{1}{2}$ .

The following is a trivial computation, but useful for many simplifications.

**Lemma C.27.**  $c_{\alpha+1}/c_\alpha = 2\alpha + 1$ .

### C.7.1 Forward Dynamics

**Theorem C.28.** *Suppose we have the nonlinearity  $\phi = \psi_1$ . Then  $\mathbf{p}^{(l)} = \Theta((1 + \sigma_v^2 \sigma_w^2 / 2)^l)$ , with the hidden constant depending on the initial condition.*

*Proof.* We have

$$\begin{aligned} \mathbf{p} &= \frac{1}{2} \sigma_v^2 (\sigma_w^2 \mathbf{p} + \sigma_b^2) + \sigma_a^2 + \mathbf{p} \\ &= \left( \frac{1}{2} \sigma_v^2 \sigma_w^2 + 1 \right) \mathbf{p} + \frac{1}{2} (\sigma_v^2 \sigma_b^2 + \sigma_a^2). \end{aligned}$$

By the standard method of characteristic equation, we get that

$$\mathbf{p}^{(l)} = A + CB^l$$

where  $A = -\frac{\sigma_a^2 + \sigma_b^2 \sigma_v^2}{\sigma_v^2 \sigma_w^2}$ ,  $B = 1 + \frac{\sigma_v^2 \sigma_w^2}{2}$ , and  $C$  is a coefficient determined by initial conditions.  $\square$

**Theorem C.29.** *Suppose  $\alpha < 1$ . We have the following asymptotic expansion*

$$\mathbf{p}^{(l)} = K_1 l^{\frac{1}{1-\alpha}} + R(l)$$

where the remainder term

$$R(l) \sim \begin{cases} -K_2 l^{\frac{\alpha}{1-\alpha}} \log l & \text{if } \alpha > \frac{1}{2} \\ (C - K_2) l \log l & \text{if } \alpha = \frac{1}{2} \text{ and } K_2 \neq C \\ \frac{C(1-\alpha)}{1-2\alpha} l & \text{if } \alpha < \frac{1}{2} \end{cases}$$

where  $K_1 = [\sigma_v^2 \sigma_w^{2\alpha} c_\alpha (1-\alpha)]^{\frac{1}{1-\alpha}}$ ,  $K_2 = \frac{1}{2} [\sigma_v^2 c_\alpha \sigma_w^{2\alpha}]^{\frac{1}{1-\alpha}} (1-\alpha)^{\frac{\alpha}{1-\alpha}} - 1$  and  $C = \sigma_a^2$ .

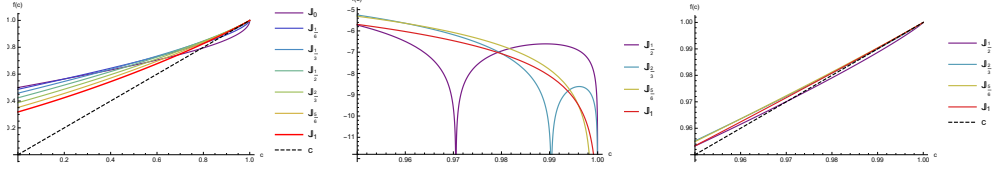
**Fig. C.14** verifies the leading coefficient and the exponent of the leading term.

*Proof.* The difference equation governing the evolution of  $\mathbf{p}$  is

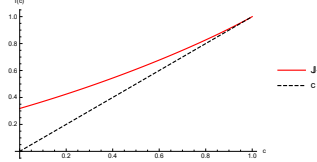
$$\mathbf{p} - \mathbf{p} = A(\mathbf{p} + B)^\alpha + C$$

where  $A = \sigma_v^2 c_\alpha \sigma_w^{2\alpha}$ ,  $B = \sigma_b^2 / \sigma_w^2$ , and  $C = \sigma_a^2$ . Then **Lemma C.15** yields the result.  $\square$

**Thm C.29** combined with **Thm C.28** gives the following result.



**Figure C.15:** (a)  $\mathbb{J}_\alpha$  for different  $\alpha$ s and the identity function. From this plot, it looks like  $\mathbb{J}_\alpha(c) \geq c$  and  $\mathbb{J}_\alpha(c) \leq 1$  for all  $\alpha \in (\frac{1}{2}, 1]$  with equality iff  $c = 1$ , but this is misleading. (b) shows  $|\mathbb{J}_\alpha(c) - c|$  in log scale. Where the curves dip below the x-axis indicate points where  $\mathbb{J}_\alpha(c) = c$ . We see that in fact every  $\mathbb{J}_\alpha$  has a solution  $\mathbb{J}_\alpha(c) = c$  for a  $c < 1$ , when  $\alpha < 1$ . (c) Furthermore, at each such  $c$ ,  $\mathbb{J}_\alpha < 1$ .



**Figure C.16:**  $\mathbb{J}_1$  vs identity

**Theorem B.16.** Suppose we have the nonlinearity  $\phi = \psi_\alpha$ . The in an FRN: If  $\alpha = 1$ , then  $\mathbf{p}^{(l)} = \Theta((1 + \sigma_v^2 \sigma_w^2 / 2)^l)$ , with the hidden constant depending on the initial condition. If  $0 < \alpha < 1$ , then  $\mathbf{p}^{(l)} = \Theta(l^{\frac{1}{1-\alpha}})$ . More precisely,  $\lim_{l \rightarrow \infty} \mathbf{p} / l^{\frac{1}{1-\alpha}} = [\sigma_v^2 \sigma_w^2 c_\alpha (1 - \alpha)]^{\frac{1}{1-\alpha}}$ .

By [2], we know that  $W\psi_\alpha(q, qc) = V\psi_\alpha(q)\mathbb{J}_\alpha(c)$ , where  $\mathbb{J}_\alpha(c) = J_\alpha(\arccos c)$  and

$$J_\alpha(\theta) := \frac{1}{2\pi c_\alpha} (\sin \theta)^{2\alpha+1} \Gamma(\alpha+1) \int_0^{\pi/2} \frac{d\eta \cos^\alpha \eta}{(1 - \cos \theta \cos \eta)^{1+\alpha}}. \quad (\Delta)$$

Note that  $\mathbb{J}_\alpha(c) \in (-\infty, \infty)$  for  $\alpha \in (-1, \infty)$  and any  $c \in (0, 1)$ , even though  $V\psi_\alpha$  is only defined for  $\alpha > -1/2$ .

**Fig. C.15** shows a comparison of  $\mathbb{J}_\alpha$  for different  $\alpha$ s along with the identity function. By [3, Lemma 11],  $\mathbb{J}_\alpha$  is an increasing and convex function as long as  $\psi_\alpha^2$  is Gaussian-integrable, which is precisely when  $\alpha > -1/2$ . We can compute  $\mathbb{J}_\alpha(1) = W\psi_\alpha(q, q)/V\psi_\alpha(q) = 1$ , and  $\mathbb{J}_\alpha(0) = W\psi_\alpha(q, 0)/V\psi_\alpha(q) = V\psi_{\alpha/2}(q)^2/V\psi_\alpha(q) = c_{\alpha/2}^2/c_\alpha = \frac{1}{2\sqrt{\pi}} \frac{\Gamma(\frac{\alpha}{2} + \frac{1}{2})^2}{\Gamma(\alpha + \frac{1}{2})}$ . We record these observations as a lemma.

**Lemma C.30.**  $\mathbb{J}_\alpha(c)$  is an increasing and convex function for each  $\alpha > -1/2$  on  $c \in [0, 1]$ .  $\mathbb{J}_\alpha(1) = 1$  and  $\mathbb{J}_\alpha(0) = \frac{1}{2\sqrt{\pi}} \frac{\Gamma(\frac{\alpha}{2} + \frac{1}{2})^2}{\Gamma(\alpha + \frac{1}{2})}$ .

For  $\alpha = 1$ , Cho and Saul [2] computed

$$\mathbb{J}_1(c) = \frac{1}{\pi} (\sqrt{1-c^2} + (\pi - \arccos(c))c).$$

**Fig. C.16** shows a plot of  $\mathbb{J}_1$  vs identity. It has derivative  $\dot{\mathbb{J}}_1(c) = 1 - \frac{1}{\pi} \arccos c$ , which shows that  $\dot{\mathbb{J}}_1(c) < 1$  with equality iff  $c = 1$ , and consequently  $\mathbb{J}_1(c) \geq c$  with equality iff  $c = 1$ . At the same time,  $\dot{\mathbb{J}}_1(c) \geq 0$  with equality iff  $c = -1$ , so  $\mathbb{J}_1$  is increasing on  $[-1, 1]$ . It has an asymptotic expansion  $\mathbb{J}_1(1 - \varepsilon) = 1 - \varepsilon + \frac{2\sqrt{2}}{3\pi} \varepsilon^{3/2} + \Theta(\varepsilon^{5/2})$  at 1.

The zeroth Bessel function of the second kind is defined by  $\mathcal{K}_0(z) = \int_1^\infty e^{-zx} (x^2 - 1)^{-1/2} dx$ . It is one of the fundamental solutions to the homogeneous differential equation  $x^2 \dot{y} + xy - x^2 y = 0$ . The following lemma shows that  $J_\alpha$  can be expressed in terms of  $\mathcal{K}_0$ .

**Lemma C.31.** For any  $\alpha > -1$ ,  $J_\alpha(\theta) = \frac{1}{2\pi c_\alpha} \sin^{2\alpha+1} \theta \int_0^\infty dx \mathcal{K}_0(x) e^{x \cos \theta} x^\alpha$

*Proof.* Cho and Saul [2] gave the expression

$$2\pi c_\alpha J_\alpha(\theta) = \csc \theta \int_0^\infty du \int_0^\infty dv e^{-(u^2+v^2-2uv \cos \theta)/2 \sin^2 \theta} u^\alpha v^\alpha.$$

Note that the integrand is symmetric in  $u$  and  $v$ . Thus, if  $V = \{(u, v) : u, v \geq 0 \ \& \ v \geq u\}$ , then

$$2\pi c_\alpha J_\alpha(\theta) = 2 \csc \theta \int_V du dv e^{-(u^2+v^2-2uv \cos \theta)/2 \sin^2 \theta} u^\alpha v^\alpha.$$

Now make the change of variables from  $V$  to  $\{(\mathbb{p}, \mathbb{q}) : \mathbb{q} \geq 2\sqrt{\mathbb{p}}\}$ :

$$\begin{aligned} \mathbb{p} &= uv & d\mathbb{p} &= v du + u dv \\ \mathbb{q} &= u + v & d\mathbb{q} &= du + dv \\ d\mathbb{p} d\mathbb{q} &= (v - u) du dv & du dv &= (\mathbb{q}^2 - 4\mathbb{p})^{-1/2} d\mathbb{p} d\mathbb{q} \end{aligned}$$

so that we have

$$2\pi c_\alpha J_\alpha(\theta) = 2 \csc \theta \int_0^\infty d\mathbb{p} e^{\mathbb{p}(1+\cos \theta) \csc^2 \theta} \mathbb{p}^\alpha \int_{2\sqrt{\mathbb{p}}}^\infty d\mathbb{q} e^{-\mathbb{q}^2 \csc^2 \theta} (\mathbb{q}^2 - 4\mathbb{p})^{-1/2}.$$

The inner integral in  $\mathbb{q}$  can be expressed in terms of  $\mathcal{K}_0$  by a change of variable  $x = \mathbb{q}^2/2\sqrt{\mathbb{p}}$ :

$$\begin{aligned} 2\pi c_\alpha J_\alpha(\theta) &= 2 \csc \theta \int_0^\infty d\mathbb{p} e^{\mathbb{p}(1+\cos \theta) \csc^2 \theta} \mathbb{p}^\alpha \frac{1}{2} e^{-\mathbb{p} \csc^2 \theta} \mathcal{K}_0(\mathbb{p} \csc^2 \theta) \\ &= \csc \theta \int_0^\infty d\mathbb{p} \mathcal{K}_0(\mathbb{p} \csc^2 \theta) e^{\mathbb{p} \cos \theta \csc^2 \theta} \mathbb{p}^\alpha \\ &= \sin^{2\alpha+1} \theta \int_0^\infty dx \mathcal{K}_0(x) e^{x \cos \theta} x^\alpha \end{aligned}$$

□

Define  $L_\alpha(\theta) = 2\pi c_\alpha J_\alpha(\theta) \csc^{2\alpha+1} \theta = \int_0^\infty dx \mathcal{K}_0(x) e^{x \cos \theta} x^\alpha$ .

**Lemma C.32.** *If  $\alpha > 1$ , then*

$$L_\alpha(\theta) = \csc^2 \theta [(2\alpha - 1) \cos \theta L_{\alpha-1}(\theta) + (\alpha - 1)^2 L_{\alpha-2}(\theta)].$$

*Proof.* We will prove this claim for  $\theta < 1$ , and by continuity this also proves the case  $\theta = 1$ . As remarked above,  $\mathcal{K}_0(z) = \ddot{\mathcal{K}}_0(z) + z^{-1} \dot{\mathcal{K}}_0(z)$ . Thus

$$\begin{aligned} L_\alpha(\theta) &= \int_0^\infty dx (\ddot{\mathcal{K}}_0(x) + x^{-1} \dot{\mathcal{K}}_0(x)) e^{x \cos \theta} x^\alpha \\ &= \dot{\mathcal{K}}_0 e^{x \cos \theta} x^\alpha \Big|_0^\infty + \mathcal{K}_0 e^{x \cos \theta} x^{\alpha-1} \Big|_0^\infty \\ &\quad - \int dx [\cos \theta e^{x \cos \theta} x^\alpha + \alpha e^{x \cos \theta} x^{\alpha-1}] \dot{\mathcal{K}}_0 \\ &\quad - \int dx [\cos \theta e^{x \cos \theta} x^{\alpha-1} + (\alpha - 1) e^{x \cos \theta} x^{\alpha-2}] \mathcal{K}_0 \end{aligned}$$

Asymptotically,  $\mathcal{K}_0(z) \sim \sqrt{\frac{\pi}{2z}} e^{-z}$  as  $z \rightarrow \infty$  and  $\mathcal{K}_0(z) \sim -\ln(z)$  as  $z \searrow 0$ , and  $\dot{\mathcal{K}}_0(z) \sim -\sqrt{\frac{\pi}{2z}} e^{-z}$  as  $z \rightarrow \infty$  and  $\dot{\mathcal{K}}_0(z) \sim -z^{-1}$  as  $z \searrow 0$ . Thus, as  $\alpha > 1$ ,

$$\begin{aligned} \dot{\mathcal{K}}_0 e^{x \cos \theta} x^\alpha \Big|_0^\infty &= - \lim_{x \rightarrow \infty} \sqrt{\pi/2} e^{-x(1-\cos \theta)} x^{\alpha-1} + \lim_{x \searrow 0} e^{x \cos \theta} x^{\alpha-1} = 0 \\ \mathcal{K}_0 e^{x \cos \theta} x^{\alpha-1} \Big|_0^\infty &= - \lim_{x \rightarrow \infty} \sqrt{\pi/2} e^{-x(1-\cos \theta)} x^{\alpha-2} + \lim_{x \searrow 0} e^{x \cos \theta} x^{\alpha-1} \ln x = 0 \end{aligned}$$

So

$$L_\alpha(\theta) = -\cos \theta L_{\alpha-1}(\theta) - (\alpha - 1) L_{\alpha-2}(\theta) - \int dx [\cos \theta e^{x \cos \theta} x^\alpha + \alpha e^{x \cos \theta} x^{\alpha-1}] \dot{\mathcal{K}}_0$$

Via another integration by parts, the integral on the right is

$$\begin{aligned} &\cos \theta e^{x \cos \theta} x^\alpha \mathcal{K}_0 \Big|_0^\infty + \alpha e^{x \cos \theta} x^{\alpha-1} \mathcal{K}_0 \Big|_0^\infty \\ &\quad - \int dx [\cos^2 \theta e^{x \cos \theta} x^\alpha + 2\alpha \cos \theta e^{x \cos \theta} x^{\alpha-1} + \alpha(\alpha - 1) e^{x \cos \theta} x^{\alpha-2}] \mathcal{K}_0 \\ &= -[\cos^2 \theta L_\alpha(\theta) + 2\alpha \cos \theta L_{\alpha-1}(\theta) + \alpha(\alpha - 1) L_{\alpha-2}(\theta)] \end{aligned}$$

where the evaluation terms vanish just like before. Altogether, we have

$$\begin{aligned} L_\alpha(\theta) &= \cos^2 \theta L_\alpha(\theta) + (2\alpha - 1) \cos \theta L_{\alpha-1}(\theta) + (\alpha - 1)^2 L_{\alpha-2}(\theta) \\ &= \csc^2 \theta [(2\alpha - 1) \cos \theta L_{\alpha-1}(\theta) + (\alpha - 1)^2 L_{\alpha-2}(\theta)] \end{aligned}$$

□

As a corollary we get

**Lemma C.33.** *Suppose  $\alpha > 1$ . Then*

$$\begin{aligned} J_\alpha(\theta) &= \cos \theta J_{\alpha-1}(\theta) + (\alpha - 1)^2 (2\alpha - 1)^{-1} (2\alpha - 3)^{-1} \sin^2 \theta J_{\alpha-2}(\theta) \\ \mathbb{J}_\alpha(c) &= c \mathbb{J}_{\alpha-1}(c) + (\alpha - 1)^2 (2\alpha - 1)^{-1} (2\alpha - 3)^{-1} (1 - c^2) \mathbb{J}_{\alpha-2}(c) \end{aligned}$$

The derivative of  $J_\alpha(\theta)$  turns out to be quite simple.

**Lemma C.34.** *Suppose  $\alpha > 0$ . Then*

$$\begin{aligned} \dot{J}_\alpha(\theta) &= -\alpha^2 (2\alpha - 1)^{-1} J_{\alpha-1}(\theta) \sin \theta \\ \dot{\mathbb{J}}_\alpha(c) &= \alpha^2 (2\alpha - 1)^{-1} \mathbb{J}_{\alpha-1}(c) \end{aligned}$$

*Proof.* We will prove the first formula. The second follows from chain rule. By [Lemma C.31](#),

$$\begin{aligned} J_\alpha(\theta) &= \frac{1}{2\pi c_\alpha} \sin^{2\alpha+1} \theta \int dx \mathcal{K}_0(x) e^{x \cos \theta} x^\alpha \\ \dot{J}_\alpha(\theta) &= \frac{1}{2\pi c_\alpha} [(2\alpha + 1) \sin^{2\alpha} \theta \cos \theta \int dx \mathcal{K}_0(x) e^{x \cos \theta} x^\alpha \\ &\quad - \sin^{2\alpha+2} \theta \int dx \mathcal{K}_0(x) e^{x \cos \theta} x^{\alpha+1}] \\ &= (2\alpha + 1) \cot \theta J_\alpha(\theta) - \frac{c_{\alpha+1}}{c_\alpha} \csc \theta J_{\alpha+1}(\theta) \\ &= (2\alpha + 1) \csc \theta [\cos \theta J_\alpha(\theta) - J_{\alpha+1}(\theta)]. \end{aligned}$$

As  $\alpha + 1 > 1$ , by [Lemma C.33](#), this is

$$\begin{aligned} &-(2\alpha + 1) \csc \theta [(\alpha - 1)^2 (2\alpha + 1)^{-1} (2\alpha - 1)^{-1} \sin^2 \theta J_{\alpha-1}(\theta)] \\ &= -(\alpha - 1)^2 (2\alpha - 1)^{-1} \sin \theta J_{\alpha-1}(\theta). \end{aligned}$$

□

Thus  $\dot{\mathbb{J}}_\alpha(1) = \alpha^2 (2\alpha - 1)^{-1} \mathbb{J}_{\alpha-1}(1) = \alpha^2 (2\alpha - 1)^{-1}$  for any  $\alpha > 0$  by [Lemma C.30](#). For  $1/2 < \alpha \leq 1$ ,  $\dot{\mathbb{J}}_\alpha(1) \geq 1$  with equality iff  $\alpha = 1$ , and for  $\alpha = 1/2$ ,  $\dot{\mathbb{J}}_\alpha(1) = \infty > 1$  by continuity of  $\dot{\mathbb{J}}_\alpha(c)$  in  $\alpha$ . Because for  $\alpha > -1/2$ ,  $\mathbb{J}_\alpha$  is increasing and convex on  $[0, 1]$  and  $\mathbb{J}_\alpha(0) > 0$  by [Lemma C.30](#),  $\mathbb{J}_\alpha$  intersects identity at a unique point away from 1 when  $\alpha \in [1/2, 1)$ . We record this as a theorem.

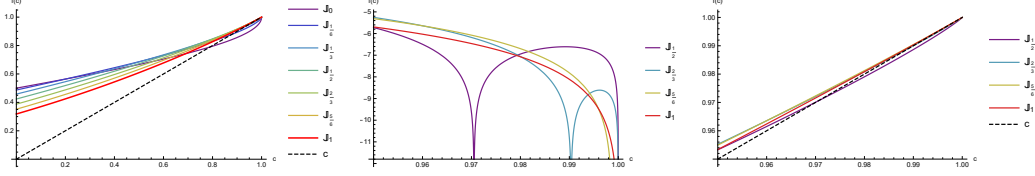
**Theorem C.35.** *For  $\alpha \in [1/2, 1)$ ,  $\mathbb{J}_\alpha(c) = c$  has two solutions: an unstable solution at 1 ("unstable" meaning  $\dot{\mathbb{J}}_\alpha(1) > 1$ ) and a stable solution in  $\mathbf{e}^* \in (0, 1)$  ("stable" meaning  $\dot{\mathbb{J}}_\alpha(\mathbf{e}^*) < 1$ ).*

This result confirms that pictures presented in [Fig. C.17b,c](#) are qualitatively correct, that there are indeed stable fixed points of  $\mathbb{J}_\alpha$  away from 1.

**Theorem B.17.** *Suppose  $\phi = \psi_1$ . Then in an FRN,  $\mathbf{e}^{(l)} \rightarrow 1$  and  $1 - \mathbf{e}^{(l)} \sim [\frac{1}{4} \sigma_v^2 \sigma_w^2 B^{-1} U l]^{-2}$  for  $B = 1 + \sigma_v^2 \sigma_w^2 / 2$  and  $U = \frac{2\sqrt{2}}{3\pi}$ . As a result,  $\mathbf{s}^{(l)} = (1 - \mathbf{e}^{(l)}) \mathbf{p}^{(l)} = \Theta(l^{-2} \exp(\Theta(l))) = \exp(\Theta(l))$ .*

*Proof.* If  $\underline{\mathbf{e}} < 1$ , then

$$\begin{aligned} c &= \frac{\sigma_w^2 \underline{\gamma} + \sigma_b^2}{\sigma_w^2 \underline{\mathbf{p}} + \sigma_b^2} \geq \underline{\mathbf{e}} \\ &\quad \mathbb{J}_1(c) \geq \mathbb{J}_1(\underline{\mathbf{e}}) \\ \underline{\mathbf{e}} &= \frac{\sigma_v^2 c_\alpha \mathbf{q}^\alpha \mathbb{J}_1(c) + \sigma_b^2}{\sigma_v^2 c_\alpha \mathbf{q}^\alpha + \sigma_b^2} \geq \mathbb{J}_1(\underline{\mathbf{e}}) \end{aligned}$$



**Figure C.17:** Left-to-right: (a)  $\mathbb{J}_\alpha$  for different  $\alpha$ s and the identity function (black, dashed line).  $\mathbb{J}_1$  is highlighted in red. From this plot, it looks like  $\mathbb{J}_\alpha(c) \geq c$  and  $\mathbb{J}_\alpha(c) \leq 1$  for all  $\alpha \in (\frac{1}{2}, 1]$  with equality iff  $c = 1$ , but this is misleading. (b) shows  $|\mathbb{J}_\alpha(c) - c|$  in log scale. Where the curves dip below the x-axis indicate points where  $\mathbb{J}_\alpha(c) = c$ . We see that in fact every  $\mathbb{J}_\alpha$  has a solution  $\mathbb{J}_\alpha(c) = c$  for a  $c < 1$ , when  $\alpha < 1$ . (c) Furthermore, at each such  $c$ ,  $\mathbb{J}_\alpha < 1$ . (b) and (c) demonstrate the existence of stable fixed points away from 1 for  $\mathbb{J}_\alpha$ ,  $\alpha \in (1/2, 1)$ , which is confirmed rigorously by [Thm C.35](#).

but  $\underline{e} \geq \mathbb{J}_1(\underline{e}) > \underline{e}$  as noted above. Thus by monotone convergence  $\underline{e}$  converges, and  $\underline{e}^* = 1$  is the only possible fixed point.

By [Lemma C.1](#),  $c = \underline{e}(1 + \Theta(\underline{\epsilon}\underline{\mathbf{p}}^{-1})) = 1 - \underline{\epsilon} + \Theta(\underline{\epsilon}\underline{\mathbf{p}}^{-1}) = 1 - u\underline{\epsilon}$  where  $u := 1 - \Theta(\underline{\mathbf{p}}^{-1})$ . Using the asymptotic expansion  $\mathbb{J}_1(1 - \epsilon) = 1 - \epsilon + U\epsilon^{3/2} + \Theta(\epsilon^{5/2})$ , we have

$$\begin{aligned}
(1 - \epsilon)\underline{\mathbf{p}} &= \sigma_v^2 \frac{\mathbf{q}}{2} \mathbb{J}_1(1 - u\underline{\epsilon}) + \sigma_a^2 + (1 - \underline{\epsilon})\underline{\mathbf{p}} \\
-\epsilon\underline{\mathbf{p}} &= \sigma_v^2 \frac{\mathbf{q}}{2} (\mathbb{J}_1(1 - u\underline{\epsilon}) - 1) - \underline{\epsilon}\underline{\mathbf{p}} \\
&= \sigma_v^2 \frac{\mathbf{q}}{2} [-u\underline{\epsilon} + Uu^{3/2}\underline{\epsilon}^{3/2} + \Theta(u^{5/2}\underline{\epsilon}^{5/2})] - \underline{\epsilon}\underline{\mathbf{p}} \\
\epsilon &= \frac{1}{\underline{\mathbf{p}}} [\underline{\mathbf{p}} + \sigma_v^2 \frac{\mathbf{q}}{2} (u - Uu^{3/2}\underline{\epsilon}^{1/2} + \Theta(u^{5/2}\underline{\epsilon}^{3/2}))] \\
&= \frac{1}{\underline{\mathbf{p}}} [\underline{\mathbf{p}} - \sigma_a^2 + \sigma_v^2 \frac{\mathbf{q}}{2} (\Theta(\underline{\mathbf{p}}^{-1}) - Uu^{3/2}\underline{\epsilon}^{1/2} + \Theta(u^{5/2}\underline{\epsilon}^{3/2}))] \\
&= \underline{\epsilon} [1 + \frac{-\sigma_a^2 + \sigma_v^2 \frac{\mathbf{q}}{2} (\Theta(\underline{\mathbf{p}}^{-1}) - Uu^{3/2}\underline{\epsilon}^{1/2} + \Theta(u^{5/2}\underline{\epsilon}^{3/2}))}{\underline{\mathbf{p}}}] \\
&= \underline{\epsilon} [1 + \frac{-\sigma_a^2 \mathbf{q}^{-1} + \frac{1}{2} \sigma_v^2 (\Theta(\underline{\mathbf{p}}^{-1}) - Uu^{3/2}\underline{\epsilon}^{1/2} + \Theta(u^{5/2}\underline{\epsilon}^{3/2}))}{\underline{\mathbf{p}}\mathbf{q}^{-1}}]
\end{aligned}$$

Let the content of the bracket on the RHS be  $\aleph$ . We have  $\underline{\mathbf{p}}\mathbf{q}^{-1} = (1 + o(1))B/\sigma_w^2$ . If  $\epsilon = O(\underline{\mathbf{p}}^{-1})$ , then  $\aleph = 1 - O(\underline{\mathbf{p}}^{-1})$ , but because  $\underline{\mathbf{p}}$  is exponentially decreasing, this means  $\epsilon = \Theta(1)$  and does not converge to 0 — this is a contradiction. Therefore,  $\underline{\epsilon} = \omega(\underline{\mathbf{p}}^{-1})$ , and

$$\begin{aligned}
\epsilon &= \underline{\epsilon} [1 - \frac{1}{2} B^{-1} \sigma_v^2 \sigma_w^2 U \underline{\epsilon}^{1/2} (1 + o(1))] \\
\epsilon - \underline{\epsilon} &= -\frac{1}{2} B^{-1} \sigma_v^2 \sigma_w^2 U \underline{\epsilon}^{3/2} (1 + o(1))
\end{aligned}$$

Using [Lemma C.14](#) to upper and lower bound our dynamics, we get that  $\epsilon^{(l)} \sim [\frac{1}{4} \sigma_v^2 \sigma_w^2 B^{-1} U]^{-2}$ .  $\square$

**Lemma C.36.** Let  $\phi$  be any nonlinearity. Suppose  $W\phi(r, rd) = V\phi(r)\mathbb{K}(d)$  for some twice differentiable function  $\mathbb{K}(d)$  independent of  $\mathbf{q}$ , where  $\mathbb{K}(1) = 1$  naturally. Suppose further that

- $\mathbb{K}(d) = d$  has a solution  $d = \mathbf{e}^* > 0$  where  $\dot{\mathbb{K}}(\mathbf{e}^*) = \delta < 1$ ;
- $\mathbb{K}(d) > d$  for all  $d < \mathbf{e}^*$  and  $\mathbb{K}(d) < d$  for all  $1 > d > \mathbf{e}^*$ ; and
- $\mathbb{K}$  is nondecreasing.

Let  $\epsilon^{(l)} := \mathbf{e}^{(l)} - \mathbf{e}^*$  and suppose  $\mathbf{e}^{(0)} < 1$ . If  $\gamma^{(l)} \rightarrow \infty$  and  $V\phi(\mathbf{q}^{(l)}) \rightarrow \infty$ , then  $\epsilon^{(l)} \rightarrow 0$  and satisfies

$$\epsilon = \underline{\epsilon} \left( 1 - \frac{\sigma_a^2 + (1 - \delta + O(\underline{\epsilon}))\sigma_v^2 V\phi(\mathbf{q})}{\underline{\mathbf{p}}} \right) + V\phi(\mathbf{q})\Theta(\gamma^{-1}\underline{\mathbf{p}}^{-1}).$$

*Proof.* First we note that because  $\mathbf{e}^*$  is the only stable fixed point of the dynamics  $x \mapsto \mathbb{K}(x)$ , with the basin of attraction  $[0, 1)$ , we can show  $\mathbf{e}^{(l)} \rightarrow \mathbf{e}^*$  as in the proof of [Thm B.11](#) (using [Lemma C.23](#)).

Write  $V^{(l)} := V\phi(\mathbf{q}^{(l)})$ . We first show that  $\mathbf{e}^{(l)} \rightarrow \mathbf{e}^*$ . When  $l$  is large,

$$\begin{aligned} c &= \frac{\sigma_w^2 \gamma + \sigma_b^2}{\sigma_w^2 \underline{\mathbf{p}} + \sigma_b^2} = \underline{c}(1 + O(\gamma^{-1})) \\ \mathbf{e} &= \frac{\sigma_v^2 V \mathbb{K}(c) + \sigma_a^2}{\sigma_v^2 V + \sigma_a^2} = \mathbb{K}(c)(1 + O(V^{-1} \mathbb{K}(c)^{-1})). \end{aligned}$$

If  $\gamma^{(l)}$  is bounded for all  $l$ , then  $\mathbf{e} \rightarrow 0$  because  $\mathbf{p}^{(l)} \rightarrow \infty$ . Since  $\mathbb{K}(c) > 0$  for  $c \in [0, 1]$  and  $V^{(l)} \rightarrow \infty$ , we have that in the limit  $l \rightarrow \infty$ ,  $\lim_{l \rightarrow \infty} \mathbf{e} = 0 = \mathbb{K}(\lim_{l \rightarrow \infty} \mathbf{e}) = \mathbb{K}(0)$  (by the continuity of  $\mathbb{K}$ ), which is impossible by our assumptions. Thus  $\gamma^{(l)} \rightarrow \infty$ , and we have  $\lim_{l \rightarrow \infty} \mathbf{e} = \mathbb{K}(\lim_{l \rightarrow \infty} \mathbf{e})$ . By our assumptions,  $\mathbf{e}^*$  is the only stable fixed point of  $\mathbb{K}$  with basin of attraction  $[0, 1)$ , so this shows that  $\mathbf{e} \rightarrow \mathbf{e}^*$  as desired.

Now we derive the equation in question. Note that  $c = \underline{c}(1 + \Theta(\gamma^{-1}))$  because  $\mathbf{e}^* < 1$ . We use the Taylor expansion  $\mathbb{K}(\mathbf{e}^* + \epsilon) = \mathbf{e}^* + \delta\epsilon + O(\epsilon^2)$ .

$$\begin{aligned} (\mathbf{e}^* + \epsilon)\mathbf{p} &= \sigma_v^2 V \mathbb{K}((\mathbf{e}^* + \epsilon)(1 + \Theta(\gamma^{-1}))) + \sigma_a^2 + (\mathbf{e}^* + \epsilon)\underline{\mathbf{p}} \\ &= \sigma_v^2 V(\mathbf{e}^* + \delta(\epsilon + \Theta(\gamma^{-1})) + O(\epsilon^2)) + \sigma_a^2 + (\mathbf{e}^* + \epsilon)\underline{\mathbf{p}} \\ \epsilon\mathbf{p} &= \sigma_v^2 V(\delta(\epsilon + \Theta(\gamma^{-1})) + O(\epsilon^2)) + \epsilon\underline{\mathbf{p}} \\ \epsilon &= \underline{c}(1 - \frac{\sigma_a^2 + (1 - \delta + O(\epsilon))\sigma_v^2 V}{\mathbf{p}}) + \Theta(V\gamma^{-1}\mathbf{p}^{-1}) \end{aligned}$$

□

**Theorem B.18.** *Suppose  $\phi = \psi_\alpha$  for  $0 < \alpha < 1$  in an FRN. Then  $\mathbf{e}$  converges to the unique nonunit fixed point  $\mathbf{e}^*$  of  $\mathbb{J}_\alpha$ , and  $|\mathbf{e}^* - \mathbf{e}^{(l)}|$  is  $\check{\Theta}(l^{-\mu})$ , where  $\mu = (1 - \dot{\mathbb{J}}_\alpha(\mathbf{e}^*)) / (1 - \alpha)$ . Additionally,  $s^{(l)} = \Theta(\mathbf{p}^{(l)}) = \Theta(l^{1/(1-\alpha)})$ .*

*Proof.* We apply [Lemma C.36](#). We first check the conditions of the lemma, with  $\mathbb{K} = \mathbb{J}_\alpha$ . The following conditions were already verified.

- $\mathbb{J}_\alpha$  has a fixed point  $\mathbf{e}^*$  less than but very close to 1, where its slope is  $v := \dot{\mathbb{J}}_\alpha(\mathbf{e}^*) < 1$ . ([Thm C.35](#))
- $\mathbb{J}_\alpha(d) > d$  for all  $d < \mathbf{e}^*$  and  $\mathbb{J}_\alpha(d) < d$  for all  $d > \mathbf{e}^*$ . (By the convexity shown in [Lemma C.30](#))
- $\mathbb{J}_\alpha$  is nondecreasing ([Lemma C.30](#)). Furthermore, from its integral formula ([Eq. \(\Delta\)](#)), we see easily that  $\mathbb{J}_\alpha$  is smooth at  $\mathbf{e}^* < 1$ .

We also proved the following

- $\mathbf{p}^{(l)} \sim [\sigma_v^2 \sigma_w^{2\alpha} c_\alpha (1 - \alpha)]^{\frac{1}{1-\alpha}} l^{\frac{1}{1-\alpha}}$  ([Thm C.29](#)) and  $\gamma^{(l)}$  is asymptotically a constant fraction of  $\mathbf{p}^{(l)}$  ([Lemma C.36](#)), so both go to  $\infty$ .
- $V\psi_\alpha(\mathbf{q}) = c_\alpha \mathbf{q}^\alpha = c_\alpha (\sigma_w^2 \mathbf{p} + \sigma_b^2)^\alpha = \Theta(l^{\alpha/(1-\alpha)})$ , so goes to  $\infty$ . ([Lemma B.15](#))

Thus, for  $v = \dot{\mathbb{J}}_\alpha(\mathbf{e}^*)$ ,

$$\begin{aligned} \frac{\sigma_a^2 + (1 - v + O(\epsilon))\sigma_v^2 V\phi(\mathbf{q})}{\mathbf{p}} &\sim \frac{(1 - v)\sigma_v^2 \sigma_w^{2\alpha} c_\alpha}{\mathbf{p}^{1-\alpha}} \\ &= l^{-1}(1 - v)/(1 - \alpha). \end{aligned}$$

Now,  $V\phi(\mathbf{q})\gamma^{-1}\mathbf{p}^{-1} = \Theta(l^{-\frac{1}{1-\alpha}-1})$ . By using the dynamics of [Lemma C.11](#) to upper and lower bound our dynamics, we have  $\epsilon^{(l)} = \Omega(l^{-\mu-\epsilon}), O(l^{-\mu+\epsilon})$  for any  $\epsilon > 0$ , where  $\mu = \min((1-v)/(1-\alpha), 1/(1-\alpha)) = (1-v)/(1-\alpha)$ . □

### C.7.2 Backward Dynamics

**Lemma C.37.** *Suppose random variable  $X \sim \mathcal{N}(0, \sigma^2)$ , and  $Y = \psi_{-\beta}(X)$  for some  $\beta > 0$ , where  $\psi_\alpha$  is  $\alpha$ -ReLU. Then for  $\xi > 0$ ,  $Y$  has density*

$$\Pr[Y \in [\xi, \xi + d\xi]] = \frac{1}{\beta\sqrt{2\pi\sigma^2}} \xi^{-\frac{1}{\beta}-1} e^{-\xi^{-2/\beta}/2\sigma^2}.$$

At  $\xi = 0$ ,  $Y$  has density given by a Dirac delta of mass  $\frac{1}{2}$ .

Furthermore,  $Y$  has finite second moment iff  $\beta < \frac{1}{2}$ .

*Proof.* We have

$$\begin{aligned} \Pr[Y \in [\xi, \infty)] &= \Pr[X \in [0, \xi^{-1/\beta}]] \\ &= \frac{1}{\sqrt{2\pi\sigma^2}} \int_0^{\xi^{-1/\beta}} e^{-x^2/2\sigma^2} dx. \end{aligned}$$

Differentiating the RHS against  $\xi$  using Leibniz's rule, we get

$$\begin{aligned} d\Pr[Y \in [\xi, \infty)]/d\xi &= \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\xi^{-2/\beta}/2\sigma^2} \frac{d}{d\xi} \xi^{-1/\beta} \\ &= \frac{-1}{\beta\sqrt{2\pi\sigma^2}} \xi^{-\frac{1}{\beta}-1} e^{-\xi^{-2/\beta}/2\sigma^2}. \end{aligned}$$

Negating both sides gives the density  $f_Y$  of  $Y$  for  $\xi > 0$ . For  $\xi = 0$ , observe that  $\lim_{\xi \rightarrow 0} f_Y(\xi) = 0$  because, while  $\xi^{-\frac{1}{\beta}-1}$  blows up polynomially,  $e^{-\xi^{-2/\beta}/2\sigma^2}$  blows up exponentially. Thus the contribution of  $Y$ 's mass at  $Y = 0$  from  $X > 0$  is 0. On the other hand, all  $X < 0$  gets mapped to  $Y = 0$ , so  $f_Y(0) = \frac{1}{2}\delta_0$ , where  $\delta_0$  is the Dirac delta.

For the second assertion, observe that

$$f_Y(\xi) \sim \frac{1}{\beta\sqrt{2\pi\sigma^2}} \xi^{-\frac{1}{\beta}-1} \text{ as } \xi \rightarrow \infty.$$

Thus,  $\xi^2 f_Y(\xi)$  is integrable iff  $2 - \frac{1}{\beta} - 1 < -1 \iff \beta < \frac{1}{2}$ . □

**Theorem B.19.** *Suppose we have the nonlinearity  $\psi_\alpha$  in an FRN.  $\text{Var}(\psi_\alpha(\zeta)^2)$  diverges for any Gaussian variable  $\zeta$  with mean 0 if  $\alpha \leq \frac{3}{4}$  but is finite if  $\alpha > \frac{3}{4}$ .*

*Proof.* Note that  $\psi_\alpha \propto \psi_{\alpha-1}$ , so it suffices to show that  $\text{Var}(\psi_{\alpha-1}(\zeta)^2) = \text{Var}(\psi_{2\alpha-2}(\zeta))$  is infinite for  $\zeta \sim \mathcal{N}(0, \sigma^2)$ . By [Lemma C.37](#) with  $\beta = 2 - 2\alpha$ ,  $\psi_{2\alpha-2}(\zeta)$  has finite variance iff  $\beta < \frac{1}{2} \iff \alpha > \frac{3}{4}$ . □

**Theorem B.20.** *Suppose we have the nonlinearity  $\psi_\alpha$  in an FRN. If  $\alpha = 1$ , then  $\chi^{(l-m)} = \chi^{(l)} (\frac{1}{2}\sigma_v^2\sigma_w^2 + 1)^m$ . If  $\alpha \in (\frac{3}{4}, 1)$ , then  $\chi^{(l-m)} = \Theta(1)\chi^{(l)}(l/(l-m))^R$  for  $R = \frac{\alpha^2}{(1-\alpha)(2\alpha-1)}$ , where the constants in  $\Theta(1)$  do not depend on  $l$  or  $m$ .*

*Proof.* If  $\alpha = 1$ , then

$$\underline{\chi} = \chi(1 + \frac{1}{2}\sigma_v^2\sigma_w^2).$$

So  $\chi^{(l-m)}/\chi^{(l)} = \Theta(1)B^m$  for  $B = 1 + \frac{1}{2}\sigma_v^2\sigma_w^2$ .



If  $\frac{1}{2} < \alpha < 1$ , then  $\underline{\chi}/\chi - 1$  is

$$\begin{aligned}
& \sigma_v^2 \sigma_w^2 V \dot{\phi}(\mathbf{q}) \\
&= \sigma_v^2 \sigma_w^2 \alpha^2 \mathbf{c}_{\alpha-1} \mathbf{q}^{\alpha-1} \\
&= \sigma_v^2 \sigma_w^2 \alpha^2 \mathbf{c}_{\alpha-1} (\sigma_w^2 \mathbf{p})^{\alpha-1} + \Theta(\mathbf{p}^{\alpha-2}) \\
&= \sigma_v^2 \sigma_w^{2\alpha} \alpha^2 \mathbf{c}_{\alpha-1} (K_1 l^{\frac{1}{1-\alpha}} - K_2 l^{\frac{\alpha}{1-\alpha}} \log l + o(l^{\frac{\alpha}{1-\alpha}} \log l))^{\alpha-1} + \Theta(l^{\frac{\alpha-2}{1-\alpha}}) \quad \text{by Thm C.29} \\
&= \sigma_v^2 \sigma_w^{2\alpha} \alpha^2 \mathbf{c}_{\alpha-1} [K_1^{\alpha-1} l^{-1} + \Theta(l^{-2} \log l)] + O(l^{-3}) \\
&= \sigma_v^2 \sigma_w^{2\alpha} \alpha^2 \mathbf{c}_{\alpha-1} K_1^{\alpha-1} l^{-1} + \Theta(l^{-2} \log l) \\
&= R l^{-1} + \Theta(l^{-2} \log l)
\end{aligned}$$

where  $R = \sigma_v^2 \sigma_w^{2\alpha} \alpha^2 \mathbf{c}_{\alpha-1} K_1^{\alpha-1} = \frac{\alpha^2}{(1-\alpha)(2\alpha-1)}$  and  $K_1 = [\sigma_v^2 \sigma_w^{2\alpha} \mathbf{c}_{\alpha-1} (1-\alpha)]^{\frac{1}{1-\alpha}}$ . So

$$\begin{aligned}
\underline{\chi} &= \chi \exp(Rl^{-1} + \Theta(l^{-2} \log l)) \\
\chi^{(l-m)} &= \Theta(1) \chi^{(l)} \left( \frac{l}{l-m} \right)^R
\end{aligned}$$

as desired. □

**Theorem B.21.** *If  $\phi = \psi_1$  in an FRN, then for  $l \geq m \geq 0$ ,*

$$\begin{aligned}
\chi_b^{(l-m)} &= \Theta(1) \chi^{(l)} B^m, & \chi_w^{(l-m)} &= \Theta(1) \chi^{(l)} B^l, \\
\chi_v^{(l-m)} &= \Theta(1) \chi^{(l)} B^l, & \chi_a^{(l-m)} &= \Theta(1) \chi^{(l)} B^m.
\end{aligned}$$

where  $B = 1 + \sigma_v^2 \sigma_w^2 / 2$ .

*If  $\phi = \psi_\alpha$  in an FRN, for  $\alpha < 1$ , then for  $l \geq m \geq 0$ ,*

$$\begin{aligned}
\chi_b^{(l-m)} &= \Theta(1) \chi^{(l)} l^R (l-m)^{-R-1}, & \chi_w^{(l-m)} &= \Theta(1) \chi^{(l)} l^R (l-m)^{\frac{\alpha}{1-\alpha}-R}, \\
\chi_v^{(l-m)} &= \Theta(1) \chi^{(l)} l^R (l-m)^{\frac{\alpha}{1-\alpha}-R}, & \chi_a^{(l-m)} &= \Theta(1) \chi^{(l)} (l/(l-m))^R.
\end{aligned}$$

*Proof.* The proof is similar to that of **Thm B.7**. □