

# SETOL: A Semi-Empirical Theory of (Deep) Learning

Charles H. Martin\*      Christopher Hinrichs†      Michael W. Mahoney‡

## Abstract

We present a Semi-Empirical Theory of Learning (SETOL) that explains the remarkable performance of State-of-the-Art (SOTA) Neural Networks (NNs). We provide a formal explanation of the origin of the fundamental quantities in the phenomenological theory of Heavy-Tailed Self-Regularization (HTSR), the Heavy-Tailed Power Law Layer Quality metrics, **AlphaHat** ( $\alpha$ ) and **AlphaHat** ( $\hat{\alpha}$ ). In prior work, these metrics have been shown to predict trends in the test accuracies of pretrained SOTA NN models, and, importantly, without needing access to the testing or even training data. Our SETOL uses techniques from Statistical Mechanics (**StatMech**) as well as advanced methods from Random Matrix Theory (RMT). Our derivation suggests new mathematical preconditions for *Ideal* learning, including the new **TRACE-LOG** metric (which is equivalent to applying the Wilson Exact Renormalization Group). We test the assumptions and predictions of our SETOL on a simple 3-layer Multi-Layer Perceptron (MLP), demonstrating excellent agreement with the key theoretical assumptions. For SOTA NN models, we show how to estimate the Model Quality of a trained NN by simply computing the Empirical Spectral Density (ESD) of the layer weight matrices and then plugging this ESD into our SETOL formulae. Notably, we examine the performance of the HTSR  $\alpha$  and the SETOL **TRACE-LOG** Layer Quality metrics, and find that they align remarkably well, both on our MLP and SOTA NNs.

---

\*Calculation Consulting, 8 Locksley Ave, 6B, San Francisco, CA 94122, [charles@CalculationConsulting.com](mailto:charles@CalculationConsulting.com).

†[chinrichs@gmail.com](mailto:chinrichs@gmail.com)

‡ICSI, LBNL, and Department of Statistics, University of California at Berkeley, Berkeley, CA 94720, [mmahoney@stat.berkeley.edu](mailto:mmahoney@stat.berkeley.edu).

# Contents

21	<b>Contents</b>	
22	<b>1 Introduction</b>	<b>5</b>
23	1.1 Statistical Mechanics (StatMech) vs. Statistical Learning Theory (SLT)	5
24	1.2 Heavy-Tailed Self-Regularization (HTSR)	6
25	1.3 What is a Semi-Empirical Theory?	7
26	1.4 A Semi-Empirical Theory of Learning (SETOL)	8
27	<b>2 Heavy-Tailed Self-Regularization (HTSR)</b>	<b>11</b>
28	2.1 The HTSR Setup	11
29	2.2 Gaussian and Heavy-Tailed Universality	13
30	2.2.1 Random Matrix Theory (RMT): Marchenko-Pastur (MP) Theory and Tracy-	
31	Widom (TW) Fluctuations	13
32	2.2.2 Heavy-Tailed Random Matrix Theory (HTRMT) and Power Law (PL) fits	14
33	2.3 Data-Free <i>Shape</i> and <i>Scale</i> Quality Metrics	16
34	<b>3 A Semi-Empirical Theory of (Deep) Learning (SETOL)</b>	<b>19</b>
35	3.1 SETOL Overview	19
36	3.2 Comparing SETOL with HTSR: Conditions for Ideal Learning	22
37	3.3 Detecting Non-Ideal Learning Conditions	23
38	3.3.1 Correlation Traps	24
39	3.3.2 Over-Regularization	25
40	<b>4 Statistical Mechanics of Generalization (SMOG)</b>	<b>28</b>
41	4.1 StatMech: the SMOG approach and the SETOL approach	28
42	4.2 Mathematical Preliminaries of Statistical Mechanics	30
43	4.2.1 Setup	31
44	4.2.2 BraKets, Expected Values, and Thermal Averages	33
45	4.2.3 Free Energies and Generating Functions	37
46	4.2.4 The Annealed Approximation (AA) and the High-Temperature Approxima-	
47	tion (high-T)	37
48	4.2.5 Average Training and Generalization Errors and their Generating Functions	40
49	4.2.6 The Quality ( $\bar{Q}$ ) and its Generating Function ( $\Gamma_{\bar{Q}}$ )	41
50	4.2.7 The Thermodynamic Large-N limit and the Saddle Point Approximation	
51	(SPA)	42
52	4.2.8 HCIZ Integrals	44
53	4.3 Student -Teacher Model	45
54	4.3.1 Operational Setup	46
55	4.3.2 Theoretical Student-Teacher Average Generalization Error ( $\bar{\mathcal{E}}_{gen}^{ST}$ ),	50
56	<b>5 Semi-Empirical Theory of the HTSR Phenomenology</b>	<b>54</b>
57	5.1 Multi-Layer Setup: MLP3	56
58	5.1.1 Data-Dependent Multi-Layer ST Self-Overlap ( $\eta(\mathbf{S}, \mathbf{T})$ )	56
59	5.1.2 A Single Layer Matrix Model	57
60	5.1.3 The Matrix-Generalized ST Overlap ( $\eta(\mathbf{S}, \mathbf{T})$ )	57
61	5.2 Quality Metrics of an Individual Layer as an HCIZ Integral	57
62	5.2.1 A Generating Function Approach to Average Quality-Squared of a Layer	57
63	5.2.2 Evaluating the Average Quality (Squared) Generating Function	59
64	5.2.3 The Effective Correlation Space (ECS)	60
65	5.2.4 Two Simplifying Assumptions: the IFA and TRACE-LOG Condition	61

66	5.3	Evaluating the Layer Quality ( $\bar{Q}$ ) in the Large- $N$ Limit	62
67	5.4	Modeling the R-Transform	64
68	5.4.1	Elementary Random Matrix Theory	64
69	5.4.2	Known R-transforms and Analytic (Formal) Models	65
70	5.4.3	Discrete Model: Spikes	65
71	5.4.4	Inverse-Wishart Model of Ideal Learning	66
72	5.4.5	Levy-Wigner Models and the AlphaHat Metric	68
73	<b>6</b>	<b>Empirical Studies</b>	<b>70</b>
74	6.1	HTSR Phenomenology: Predicting Model Quality via the Alpha metric	71
75	6.2	Testing the Effective Correlation Space	73
76	6.2.1	Train and test errors by epochs	74
77	6.2.2	Truncation and Generalization	76
78	6.3	Evaluating the Trace-Log Condition	77
79	6.3.1	The MLP3 model	77
80	6.3.2	State-of-the-Art (SOTA) models	78
81	6.4	Inducing a Correlation Trap	80
82	6.5	Overloading and the Hysteresis Effect	80
83	6.5.1	Baseline: Loading onto both FC1 and FC2	81
84	6.5.2	Overloading FC1	81
85	6.5.3	Overloading FC2	81
86	<b>7</b>	<b>Conclusion and Future Directions</b>	<b>87</b>
87	7.1	Future Directions	88
88	<b>A</b>	<b>Appendix</b>	<b>91</b>
89	A.1	Data Vectors, Weight Matrices, and Other Symbols	91
90	A.2	Summary of the Statistical Mechanics of Generalization (SMOG)	94
91	A.2.1	Annealed Hamiltonian $H^{an}(R)$ when Student and Teachers are Vectors	94
92	A.2.2	Annealed Hamiltonian $H^{an}(\mathbf{R})$ for the Single Layer Matrix-Generalized ST	
93		Error	96
94	A.3	Expressing the Layer Quality	100
95	A.4	Derivation of the TRACE-LOG Condition	101
96	A.4.1	Setting up the Saddle Point Approximation (SPA)	101
97	A.4.2	Casting the Generating Function ( $\beta \mathbf{T}_{Q^2}^{IZ}$ ) as an HCIZ Integral	104
98	A.5	MLP3 Model Details	104
99	A.6	Tanaka's Result	105
100	A.6.1	Setup and Outline	106
101	A.6.2	Step 1. Forming the Integral Transformation of ESD ( $\rho_{\mathbf{A}}^{\infty}(\lambda)$ )	107
102	A.6.3	Step 2: The Saddle Point Approximation (SPA): Explicitly forming the	
103		Large Deviation Principle (LDP)	110
104	A.6.4	Expressing the Norm Generating Function ( $\mathbb{G}_{\mathbf{A}}(\lambda)$ ) as the Integrated R-	
105		transform ( $R(z)$ ) of the Correlation Matrix ( $\mathbf{A}$ )	112
106	A.6.5	Selecting $\mathbf{A} := \mathbf{A}_1$ instead of $\mathbf{A}_2$	114
107	A.7	The Inverse-Wishart (IW) Model	114
108	A.7.1	The Branch Cut in the IW Model	114
109	A.7.2	$R(z)[IW]$ is Complex Along the Branch Cut	115
110	A.7.3	Calculation of $G(\lambda)[IW]$	115
111	A.7.4	Computing the Modulus $ G(\lambda)[IW] $	116

DRAFT

# 1 Introduction

Deep Neural Networks (DNNs)—models in the field of Artificial Intelligence (AI)—have driven remarkable advances in multiple fields of science and engineering. AlphaFold has made significant progress in solving the protein folding problem.[?] Notably, the 2024 Nobel Prize in Physics was awarded to Hopfield and Hinton for developing early approaches to AI using techniques from *Statistical Mechanics* (StatMech), and Jumper and Hassabis, along with Baker, received the 2024 Nobel Prize in Chemistry for their contributions to AlphaFold and computational protein design.[?, ?] Self-driving cars now roam the streets of major metropolitan cities like San Francisco. Large Language Models (LLMs) like ChatGPT have gained worldwide attention and initiated serious conversations about the possibility of creating an Artificial General Intelligence (AGI). Clearly, not a single area of science or engineering has ignored these remarkable advances in the field of AI and Neural Networks (NNs).

Despite this remarkable progress in a research field spanning over 50 years, developing, training, and maintaining such complex models require staggering capital resources, limiting their development to only the largest and best-funded organizations. While many such entities have open-sourced some of their largest models (such as Llama and Falcon), using these models requires assuming they have been trained optimally, without significant defects that could limit, skew, or even invalidate their use downstream. Moreover, testing such models can be very expensive and complex to interpret.

Because training and evaluating NNs is so hard, significant issues can manifest in many obvious and non-obvious ways. A primary research goal is to improve the efficiency and reduce the cost of training large NNs. A less known but critical issue arises in many industrial settings, specifically “selecting the best models to test.” This arises in industries such as ad-click prediction, search relevance, quantitative trading, and more. Frequently, one has several seemingly equally good models to choose from, but testing the model can be very expensive, time-consuming, and even risky to the business. Recently, researchers and practitioners have started to fine-tune such large open-source models using techniques such as LoRA and QLoRA. Such methods allow one to adapt a large, open-source NN to a small dataset, and very cheaply. However, in fine-tuning, one could unwittingly overfit the model to the small dataset, degrading its performance for its intended use. Despite these and many other problems, theory remains well behind practice, and there is an increasingly pressing need to develop *practical predictive theory* to both improve the training of these very large NN models and to design new methods to make their use more reliable.

Before discussing these methods, however, let us explain *What is a SemiEmpirical Theory*

## 1.1 Statistical Mechanics (StatMech) vs. Statistical Learning Theory (SLT)

Historically, there have been two competing theoretical frameworks for understanding NNs: *Statistical Mechanics* (StatMech) [?, ?, ?, ?, ?, ?, ?]; and *Statistical Learning Theory* (SLT) [?].

- **Statistical Mechanics (StatMech).** This framework has been foundational to the early development of NN models, such as the Hopfield Associative Memory (HAM) [?], Boltzmann Machines [?], [?], etc. StatMech has also been used to build early theories of learning, such as the Student-Teacher model for the Perceptron Generalization Error [?, ?], the Gardner model [?], and many others. Notably, the HAM was based on an idea by Little, who observed that, in a simple model, long-term memories are stored in the eigenvectors of transfer matrix [?]. (This general idea, but in a broader sense, is central to our approach below.) Moreover, StatMech predicts that NNs exhibit phase behavior. This has recently been rediscovered as the Double Descent phenomenon [?, ?], but it was known in StatMech long before it’s recent rediscovery [?]. However, unlike other applied physics theories (e.g.,

Semi-Empirical methods in quantum chemistry), **StatMech** only offers qualitative analogies, failing to provide lacks quantitative predictions about large, modern NN models.[?]

- **Statistical Learning Theory (SLT)**. SLT and related approaches (VC theory, PAC bounds theory, etc.) have been developed within the context of traditional computational learning problems [?], and they are based on analyzing the convergence of frequencies to probabilities (over problem classes, etc.). It was recognized early on, however, that they could not be directly applied to NNs [?]. Moreover, SLT cannot even reproduce quantitative properties of learning curves [?, ?] (whereas **StatMech** is very successful at this [?]). SLT also failed to predict the “Double Descent” phenomena [?]. More recently, it has been shown that in practical settings SLT can give vacuous [?] or even opposite results to those actually observed [?].

Technically, SLT focuses on obtaining bounds on a model’s worst-case behavior, while **StatMech** seeks a probabilistic understanding of typical behaviors across different states or configurations. Unfortunately, neither of these general theoretical frameworks has proven particularly useful to NN practitioners. **SETOL** combines insights from both. Rather than being purely phenomenological like the **HTSR** approach, **SETOL** is derived from first-principles, and in form of a Semi-Empirical theory. As such, **SETOL** offers a practical, Semi-Empirical framework that bridges rigorous theoretical modeling and empirical observations for modern NNs.

## 1.2 Heavy-Tailed Self-Regularization (HTSR)

**HTSR** theory is an approach that combines ideas from **StatMech** with those of *Heavy-Tailed Random Matrix Theory* (RMT), providing eigenvalue-based quality metrics that correlate with model quality (i.e., out-of-sample performance). **HTSR** theory posits that well-trained models have extracted subtle correlations from the training data, and that these correlations manifest themselves in the *Shape* and *Scale* of the eigenvalues of the layer weight matrices **W**. In particular, if one computes the empirical distribution of the eigenvalues,  $\lambda_i$ , of an individual  $N \times M$  weight matrix, **W**, then this density,  $\rho^{emp}(\lambda)$ , which is an ESD, is Heavy-Tailed (HT) and can be well-fit to a *Power Law* (PL), i.e.,  $\rho(\lambda) \sim \lambda^{-\alpha}$ , with exponent  $\alpha$ . **HTSR** theory provides a *phenomenology* for qualitatively-distinct phases of learning [?]. It can, however, also be used to define *Layer-level Quality metrics* and *Model-level Quality metrics*: e.g., the **Alpha** ( $\alpha$ ) and **AlphaHat** ( $\hat{\alpha}$ ) PL metrics, described below.

Not needing any training data, **HTSR** theory has many practical uses. It can be directly applicable to large, open-source models where the training and test data may not be available. Model quality metrics can be used, e.g., to predict trends in the quality of SOTA models in computer vision (CV) [?] and natural language processing (NLP) [?, ?], both during and after training, and without needing access to the model test or training data. Layer quality metrics can be used to diagnose potential internal problems in a given model, or (say) to accelerate training by providing optimal layer-wise learning rates [?] or pruning ratios [?]. Most notably, the **HTSR** theory provides *Universal Layer Quality metrics* encapsulated in what appears to be a critical exponent,  $\alpha = 2$ , that is empirically associated with optimal or Ideal Learning. Moreover, as argued below, the value  $\alpha = 2$  appears to define a phase boundary between a generalization and overfitting, analogous to the phase boundaries seen in **StatMech** theories of NN learning.

These results both motivate the search for a first principles understanding of the **HTSR** theory, and suggests a path for developing a practical predictive theory of Deep Learning. For this, however, we need to go beyond the phenomenology provided by **HTSR** theory, to relate it to some sort of (at least semi-rigorous/semi-empirical) derivations based on the **StatMech** theory of learning, and drawing

upon previous success (in Quantum Chemistry) in developing a first principles Semi-Empirical theory.

### 1.3 What is a Semi-Empirical Theory?

Historically, one of the most well known *Semi-Empirical* methods comes from Nuclear Physics. The Semi-Empirical Mass Formula, dating back to 1935, is based on the heuristic Liquid Drop Model of the nucleus, and it was used to predict experimentally-observed binding energies of nucleons. This model describes nuclear fission, and it was central to its development of the atomic bomb:

Prior to WWII, Nuclear Physics was a phenomenological science, which relied upon experimental data and descriptive models [?].

In the Post-war era, the epistemological nature of nuclear theory changed, as it saw the development of Semi-Empirical shell models of the nucleus. These models were formulated with rigor (in the physics sense) but also relied on heuristic assumptions and experimental data for accurate predictions. They captured the structure of atomic nuclei and could accurately describe various nuclear properties[?, ?, ?]. The shell models, analogous to the electronic shell structure of atoms, represented a shift toward a more rigorous understanding of nuclear phenomena.

About this time, RMT itself was also introduced by *Wigner* [?] to model the statistical patterns of the nuclear energy spectra of strongly interacting heavy nuclei. These patterns were universal, independent of the specific nucleus, suggesting that a probabilistic approach would be fruitful. In the following decades, RMT saw many advances, including the development of the Marchenko-Patur model[?], and numerous other applications in physics[?]. By the 1990s, RMT was further expanded when *Zee* introduced the *Blue Function*, and reinterpreted the *R-transform* as a self-energy within the framework of many-body / quantum field theory (QFT) [?]. Also, so-called HCIZ integrals, integrals over random matrices, were being used both to model disordered electronic spectra[?], and, later, the behavior of spin glass models[?, ?].

Returning to the 1950s, and prior to the development of highly accurate, modern, computational *ab initio* theories of Quantum Chemistry, theoretical chemists introduced the Semi-Empirical PPP method for conjugated polyenes [?]. The PPP model recast the electronic structure problem as an *Effective Hamiltonian* for the  $\pi$ -electrons.<sup>1</sup> For many years this and related Semi-Empirical methods worked remarkably well, even better than the existing *ab initio* theories[?, ?, ?, ?]. Most importantly, these methods could be *fit* on a broad set of empirical molecular data, and then applied to molecules not in the original training set.

Around the same time, *Löwdin* first formalized the concept of the Effective Hamiltonian, which allowed the reduction of complex many-body problems to simplified *Effective Potentials* that still captured the essential physics. Then in the late 1960s Brandow developed an Effective Hamiltonian theory of nuclear structure, leveraging the *Linked Cluster Theorem* (LCT) (see [?]) and quantum mechanical many-body theory to describe the highly correlated effective interactions in a reduced model space.<sup>2</sup>

Like modern NNs, these Semi-Empirical methods of Quantum Chemistry worked well beyond their apparent range of validity, generalizing very well to out-of-distribution (OOD) data. This led to the search for a Semi-Empirical *Theory* to explain the remarkable performance of these phenomenological methods. Building on Brandow’s many-body formalism, Freed and collaborators [?, ?] developed an *ab initio* Effective Hamiltonian Theory of Semi-Empirical methods to explain

<sup>1</sup>The PPP model resembles the later developed tight-binding model of condensed matter physics[?]

<sup>2</sup>Note also that the LCT shows that the log partition function (*i.e.*,  $\ln Z$ ) can be expressed a sum of connected diagrams, which is very similar to our result below, which expresses the log partition function here as a sum of matrix cumulants from RMT.

the remarkable success of the Semi-Empirical methods. Specifically, the values of the PPP empirical parameters could be directly computed effective interactions, including both renormalized self-energies and higher-order terms. Somewhat later, in the 90s, Martin et. al.[?, ?, ?, ?] extended and applied this Effective Hamiltonian theory and demonstrated the Universality of the Semi-Empirical PPP parameters numerically. Indeed, it is this Universality that enabled the for-a-time inexplicable ODD performance of these methods. Crucially, this decades long line of work established a comprehensive analytic and numerical *Theory* of Semi-Empirical methods. That is, a framework that confirmed the empirically observed Universality, provided theoretical justification for this, and enabled systematic improvements of the methods using numerical techniques.

Finally, it is important to mention the Effective Hamiltonian approach provided by the Wilson *Renormalization Group* (RG).[?] The RG approach provides a powerful framework for studying strongly correlated systems across different scales, enabling the construction of an Effective Hamiltonian by *integrating out* weakly-correlated degrees of freedom in a Scale-Invariant way. It is particularly suited for critical points and phase boundaries—such the phase boundary between generalization and memorization in spin glass models of neural networks— and, importantly, predicts the existence of Universal Power Law (PL) exponents .

**Relevance to Deep Learning** In this sense, Semi-Empirical theories of Nuclear Physics and Quantum Chemistry, (as well as the Renormalization Group approach), seem particularly appropriate for Deep Learning. DNN models are complex black boxes that defy statistical descriptions they are commonly pre-trained on a large set of data; and then applied to new data sets in new domains via transfer learning. Most recently, the inexplicable success of transfer-learning is seen in the GPT (Generative Pre-trained Transformer) models[?], and motivated early work by Jumper et. al. on protein folding[?]

In contrast, these Semi-Empirical approaches differ from more recently developed theoretical approaches to deep learning, which are typically based on SLT, rather than **StatMech** [?]. In particular, there have recently appeared several theories of deep learning, formulated using ideas from RMT. However, regarding realistic models, it has been explicitly stated that “These networks are however too complex in general for developing a rigorous theoretical analysis on the spectral behavior [?]. Even in recent work applying RMT to NNs, it has been noted “*that we make no claim about trained weights, only random weights*” [?]. The weight matrices of a trained NN, however, are clearly *not* simply random matrices—since they encode the specific correlations from the training data.

## 1.4 A Semi-Empirical Theory of Learning (SETOL)

We propose SETOL, a Semi-Empirical Theory for Deep Learning Neural Networks (NNs), as both a theoretical foundation for HTSR phenomenology and a novel framework for predicting the properties of complex NN models. This unified framework offers a deeper understanding of DNN generalization through a Semi-Empirical approach inspired by many-body physics, combined with a classic **StatMech** model for NN generalization. Specifically, SETOL combines theoretical and empirical insights to evaluate Model Quality, showing that the weightwatcher layer HTSR PL metrics (**Alpha** and **AlphaHat**) can be derived using a phenomenological Effective Hamiltonian approach. This approach expresses the HTSR Layer Quality in terms of the RMT matrix cumulants of the layer weight matrix  $\mathbf{W}$ , and is governed by a Scale-Invariant transformation equivalent to a single step of an exact Renormalization Group (RG) transformation. Here, we derive this from first principles, requiring no previous knowledge of statistical physics.

The SETOL approach unifies the HTSR principles with a broader theoretical framework for layer analysis. The HTSR theory identifies Universality (e.g.,  $\alpha = 2$ ) as a hallmark of the best-trained



DNN layers, and, here, our **SETOL** introduces the closely related *Trace Log Condition*, a Scale-Invariant or Volume-Preserving transformation that reflects an underlying *Conservation Principle*. Together, these principles form the theoretical foundation for deriving HTSR Layer Quality metrics from first principles. By leveraging techniques from **StatMech** and modern RMT, **SETOL** offers a rigorous framework to connect empirical observations with theoretical predictions, advancing our understanding of generalization in neural networks.

- **Derivation of the HTSR Layer Quality metrics Alpha and AlphaHat** The **SETOL** approach takes as input the Empirical Spectral Density (ESD) of the layers of trained NN, and derives an expression for the approximate *Average Generalization Accuracy* of a multi-layer NN. We call this approximation the *Model Quality*, denoted  $\bar{Q}^{NN}$ . This Model Quality is expressed as product of individual Layer Quality terms,  $\bar{Q}_L^{NN}$ , which themselves can then be directly related to the HTSR Power Law (PL) empirical Alpha ( $\alpha$ ) and AlphaHat ( $\hat{\alpha} = \alpha \log_{10} \lambda_{max}$ ) metrics.

In particular, the Layer Quality-Squared,  $\bar{Q}^2 \approx [\bar{Q}_L^{NN}]^2$ , is expressed the logarithm of an HCIZ integral, the Thermal Average of an Annealed Error Potential for a matrix-generalized form the Linear Student-Teacher model of classical **StatMech**. This HCIZ integral evaluates into the sum of integrated R-transforms from RMT, or, equivalently, as a sum of integrated matrix cumulants. From this, the HTSR AlphaHat metric can be derived in the special case of Ideal Learning.<sup>3</sup>

- **Discovery of a Mathematical Condition for Ideal Learning.** By Ideal Learning, we mean that the specific NN layer has optimally converged, capturing as much of the information as possible in the training data without overfitting to any part of it. In defining this, and deriving our results, we have discovered (and are proposing) a new condition for Ideal Learning, which is associated with the Universality of the HTSR theory

- **HTSR Condition for Ideal Learning.** This HTSR theory states that a NN layer is Ideal when the ESD can be well fit to a Power Law (PL) distribution, with PL exponent  $\alpha = 2$ . Importantly, this appears to be a Universal property of all well-trained NNs, independent of the training data, model architecture, and training procedure.
- **SETOL TRACE-LOG Condition for Ideal Learning.** The **SETOL** condition for Ideal Learning states that the dominant eigencomponents associated with the ESD of layer form a reduced-rank *Effective Correlation Space* (ECS) that satisfies a new kind of Conservation Principle or *Volume Preserving Transformation* such that the largest eigenvalues  $\tilde{\lambda}_i$  of the ECS satisfy the condition  $\ln \prod \tilde{\lambda}_i = \sum \ln \tilde{\lambda}_i = 0$ . This is called the *TRACE-LOG Condition*. The TRACE-LOG Condition is equivalent to the taking a single step of the Wilson Exact Renormalization (RG).

The HTSR Condition has been proposed and analyzed previously [?, ?, ?]; but the TRACE-LOG Condition is new, based on our **SETOL** theory. When these two conditions align, we propose the NN layer is in the Ideal state.

- **Experimental Validation.** We present detailed experimental results on a simple model, along with observations on large-scale pretrained NNs, to demonstrate that the HTSR conditions for ideal learning ( $\alpha = 2$ ) are experimentally aligned with the independent **SETOL** condition for ideal learning ( $\det(\tilde{\mathbf{X}}) = 1$ ). See Section. 6.3. Our primary objective here is

<sup>3</sup>The **SETOL** approach to the HTSR theory resembles in spirit the derivation of the Semi-Empirical PPP models using the Effective Hamiltonian theory, where each phenomenological parameter is associated with a renormalized effective interaction, expressed as a sum of linked diagrams or clusters.[?, ?]

not to demonstrate performance improvements on SOTA NNs—this has been previously established [?]. Instead, our aim is to **validate the theoretical assumptions** of SETOL, test the **predictions of the SETOL framework**, and examine the **new, independent learning conditions** we discovered—on a model that is sufficiently simple that we can evaluate and stress test the theory.

- **Observations on Overfit Layers ( $\alpha < 2$ ).** Being a Semi-Empirical theory, SETOL can also identify violations of its assumptions. For example, when empirical results show  $\alpha < 2$  for a single layer, the layer’s ESD falls into the HTSR Very Heavy-Tailed (VHT) Universality class. (See Section 6.5.) When this happens, the layer may be slightly overfit to the training data, resulting in **suboptimal performance** and potentially even exhibiting **hysteresis-like effects** (memory effects)—that we observe empirically. These effects indicate that overfit layers may retain memory-like behavior, affecting learning dynamics and generalization.

## 2 Heavy-Tailed Self-Regularization (HTSR)

In this section, we provide an overview of the HTSR phenomenology.<sup>4</sup> HTSR has been presented in detail previously [?, ?, ?].<sup>5</sup> Here, we provide a self-contained summary, with an emphasis on certain technical issues that will be important for our SETOL. We highlight its practical application for interpreting observed behaviors in trained weight matrices, and we distinguish the HTSR phenomenology from the analytical methods used in the SETOL approach. In Section 2.1, we summarize the basic HTSR setup and results; in Section 2.2, we summarize Gaussian (for RMT) and Heavy-Tailed (for HTRMT) Universality; and in Section 2.3, we describe *Shape* metrics and *Scale* metrics that arise from HTSR. (Here, we focus on basic methods for identifying HT correlations in the ESDs of pre-trained weight matrices; in Sections 6.1, 6.4 and 6.5, we show detailed experiments using theoretical constructs from the HTSR phenomenology.)

### 2.1 The HTSR Setup

We can write the Energy Landscape function (or NN output function) for a NN with  $L$  layers as

$$E_{NN} := h_L(\mathbf{W}_L \times h_{L-1}(\mathbf{W}_{L-1} \times (\dots) + \mathbf{b}_{L-1}) + \mathbf{b}_L) \quad (1)$$

with activation functions  $h_l(\cdot)$ , and with weight matrices and biases  $\mathbf{W}_l$  and  $\mathbf{b}_l$ .<sup>6</sup> For simplicity of exposition here (HTSR can be applied much more broadly), we ignore the structural details of the layers (dense or not, convolutions or not, residual/skip connections, etc.). We also ignore the biases  $\mathbf{b}_l$  (because they can be subsumed into the weight matrices), and we treat each layer as though it contains a single weight matrix  $\mathbf{W}_L$ . We imagine training (or fine-tuning) this model on labeled data  $\{\mathbf{x}_\mu, y_\mu\} \in \mathcal{D}$ , where  $\mathbf{x}_\mu$  is the  $\mu$ -th input vector and  $y_\mu$  is its corresponding label (e.g., for binary classification,  $y_\mu \in \{-1, 1\}$ ). We expect to use backprop via some variant of stochastic gradient descent (SGD) to minimize some loss functional,  $\mathcal{L}$  (such as  $\ell_2$ , cross-entropy, etc.):

$$\operatorname{argmin}_{\mathbf{W}_l, \mathbf{b}_l} \sum_{\mu} \mathcal{L}[E_{NN}(\mathbf{x}_\mu), y_\mu] + \Omega, \quad (2)$$

where,  $\Omega$  denotes some explicit regularizer (such as an  $\ell_1$  or  $\ell_2$  constraint on layer weight matrices) or some implicit regularization procedure (such as clipping the weight matrices or applying dropout).

Given a real-valued  $N \times M$  layer weight matrix  $\mathbf{W}$  (dropping the subscript), let  $\mathbf{X}$  be the  $M \times M$  layer *Correlation Matrix*:

$$\mathbf{X} := \frac{1}{N} \mathbf{W}^\top \mathbf{W}. \quad (3)$$

The Empirical Spectral Density (ESD) of  $\mathbf{W}$ , denoted  $\rho^{emp}(\lambda)$ , is formed from the  $M$  eigenvalues  $\lambda_j$  of  $\mathbf{X}$ :

$$\rho^{emp}(\lambda) := \sum_{j=1}^M \delta(\lambda - \lambda_j). \quad (4)$$

Given a model, we can compute the ESDs of all of its layers, as well as other metrics below, with the open-source **WeightWatcher** tool [?].<sup>7</sup>

<sup>4</sup>We may also refer to the HTSR phenomenology as the HTSR Theory; we use the term phenomenology to emphasize its empirical nature, and to distinguish it from the analytical methods used in the SETOL approach.

<sup>6</sup>The Energy Landscape function  $E_{NN}$  acts on a data instance and generates a list of energies, or un-normalized probabilities; [?]. This notation was chosen to make an analogy with Random Energy Models (REM) from spin-glass and protein folding theories [?, ?].

<sup>7</sup>For practical purposes, the **WeightWatcher** tool computes  $\rho^{emp}(\lambda)$  by forming the Singular Value Decomposition (SVD) of the layer weight matrices  $\mathbf{W}$ , computing the eigenvalues  $\lambda = \sigma^2$  from the singular values  $\sigma$ , and (when useful) smoothing them with a Kernel Density Estimator (KDE). For some calculations, such as the TRACE-LOG condition, we must also select the appropriate normalization of  $\mathbf{W}$ ,

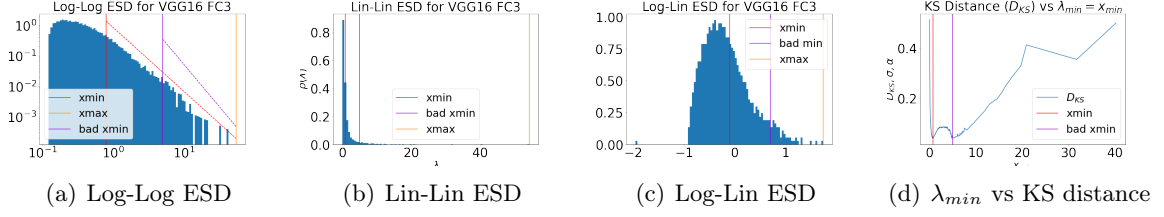


Figure 1: **Fitting ESDs within HTSR.** Depiction of the ESD and results of PL fits for a typical well-trained layer of a modern NN (FC3 of VGG19), including both the actual and good PL fit (red) and a hypothetical bad PL fit (purple). The same ESD is plotted on a Log-Log (a), Lin-Lin (b) and Log-Lin (c) scales. (d) depicts how the start of the PL tail,  $\lambda_{min}$ , varies with the quality of the PL fit (the  $D_{KS}$  distance). All plots are generated using the open-source **WeightWatcher** tool. See the main text for details.

Based on empirical results based on thousands of pre-trained models and tens of thousand of layers [?, ?, ?, ?], it is generally observed that the best performing NNs have ESDs that are HT, and whose *tails* of these ESDs,  $\rho_{tail}(\lambda)$ , can be well fit to a PL, beyond some cutoff  $\lambda \geq \lambda_{min}$ .<sup>8</sup> For a PL fit,

$$\rho_{tail}(\lambda) := \rho^{emp}(\lambda \geq \lambda_{min}) \sim \lambda^{-\alpha}, \quad (5)$$

where  $\lambda_{min}$  is where the tail of the ESD starts (i.e., it is not the minimum eigenvalue, but the minimum eigenvalue in the tail of the ESD). See Figure 1. As such, the tail of the ESD “starts” at some value  $\lambda_{min}$ , called *xmin* here, and it continues until the maximum eigenvalue  $\lambda_{max}$ , called *xmax* here (labeled *xmax* in the figure, shown by the orange line) We estimate *xmin* and  $\alpha$  jointly, using the method of [?], as implemented in the **powerlaw** python package [?], which is also integrated into the open-source **WeightWatcher** tool [?, ?].<sup>9</sup>

**Fitting ESDs.** Choosing the start of the tail,  $\lambda_{min}$ , is important for HTSR (and it will be very important for SETOL, as we will describe below). See Figure 1 for a depiction of how this was done within HTSR theory. Figures 1(a)-1(c) show the results of both a “good fit” and a “bad fit” on the same ESD, while Figure 1(d) indicates the quality of fit. For the good fit, the start of the tail is the optimal value  $\lambda_{min} = xmin$  (in red); and for the bad fit, it is a suboptimal *bad xmin* (purple). Figure 1(d) depicts how the best fit is determined; it plots  $xmin = \lambda_{min}$  versus the  $D_{KS}$  value, which is the Kolmogorov-Smirnov (KS) distance between the PL fit and the empirical data [?]. Notice that there are two nearly degenerate minima on Figure 1(d), corresponding to the good fit and the bad fit. It is not uncommon to face such practical challenges, as real-world ESDs are often slightly deformed from a perfect PL density, e.g., they may have two or more near degenerate solutions on the KS plots (d). (They may also have anomalously large eigenvalues; this is discussed in more detail in Section 3.3.)

When one finds a good PL fit for the ESD of a layer **W**, it provides information about the *Shape* and *Scale* of the ESD of that layer. In particular: the **SpectralNorm**,  $\lambda_{max}$ , being a matrix norm, is a measure of the size *Scale* of the ESD [?]; the fitted PL exponent **Alpha**,  $\alpha$ , being

<sup>8</sup>Doing a large meta-analysis like this is tricky; but see [?, ?, ?, ?]. The **WeightWatcher** tool provides a systematic, reproducible way to compute a PL fit (using an MLE method of Clauset et al. [?]), as well other model metrics, including the **SpectralNorm**, **Rand-Distance**, and **AlphaHat** metrics [?]. Also, the ESD  $\rho(\lambda)_{tail}$  is sometimes better fit by a Truncated Power Law (TPL), due to finite-size effects. This is important in practice, but we ignore this complexity in this initial discussion of SETOL.

<sup>9</sup>The authors of [?] failed to find evidence of a PL-like distribution in NN weight matrices, which is likely to be the case when  $\alpha$  and *xmin* are not estimated *jointly*, as can be seen in Figure 1(d).

HT/RMT Universality class	$\mu$ range	$\alpha$ range	Best Fit
RandomLike	NA	NA	MP
Bulk+Spikes	NA	NA	MP+Spikes
Weakly Heavy Tailed	$\mu > 4$	$\alpha > 6$	PL
Heavy (Fat) Tailed	$\mu \in (2, 4)$	$\alpha \in (2, 6)$	PL
Very Heavy Tailed	$\mu \in (0, 2)$	$\alpha \in (1, 2)$	(T)PL
Rank Collapse	NA	NA	NA

Table 1: HTSR Heavy-Tailed Universality classes of RMT. See Table 1 of [?] for more details.

the slope of the tail of the ESD on a Log-Log plot, describes the *Shape* of the ESD; and the **WeightWatcher AlphaHat** metric combines *Shape* and *Scale* information. Also, as opposed to other applications of PL fits [?, ?], in our analysis, the start of the tail,  $\lambda_{min} = \lambda_{min}^{PL}$ , plays a particularly important role because it identifies the subspace of the strongest generalizing eigenvectors (i.e.,  $\tilde{\mathbf{X}}$ , below) in each layer.

## 2.2 Gaussian and Heavy-Tailed Universality

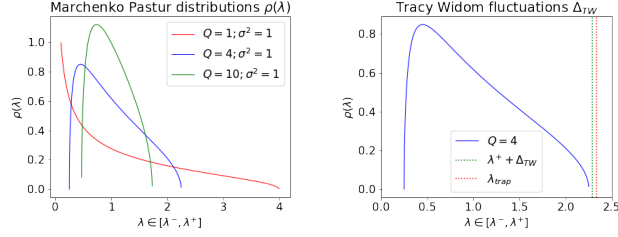
The HTSR phenomenology uses RMT to classify the ESD of a layer  $\mathbf{W}$  into one of 5+1 Phases of Training, each roughly corresponding to a (Gaussian or HT) Universality class (of RMT or HTRMT). This is summarized in Table 1. A Universality class is a set of matrices having a common limiting spectral distribution, regardless of the other properties of their entries. Of those, the most familiar is the Gaussian class, characterized by the Marchenko Pastur (MP) results from traditional RMT [?, ?]. The Gaussian Universality class, however, is particularly poorly suited for analyzing realistic NNs—precisely because the ESDs of SOTA NNs are well-fit by HT distributions. This should not be surprising: weight matrices of realistic NNs do *not* have independent (i.i.d.) entries—their entries are strongly-correlated precisely because they provide a view into the correlated training data.

To model strongly-correlated NN layer matrices, the HTSR phenomenology characterizes NN layer weight matrices in terms of their ESDs (when a good PL fit can be found) by postulating that the (tail of the) eigenvalue spectrum  $\rho(\lambda)$  determines how each layer contributes to the overall generalization. To do this, the HTSR approach models the strong-correlated layer weight matrices *as if* they are actually i.i.d. HT random (i.e., entry-wise uncorrelated) matrices. By doing this, one can associate each  $\rho(\lambda)$  with the corresponding HT Universality class, according to the PL exponent  $\alpha$  fitted from the ESD. As we will see in Section 3.3, it can be critical to distinguish when the ESD is HT *Correlation-wise* vs HT *Element-wise*.

To understand Table 1 better, we first review basic results.

### 2.2.1 Random Matrix Theory (RMT): Marchenko-Pastur (MP) Theory and Tracy-Widom (TW) Fluctuations

The Marchenko-Pastur (MP) distribution predicts the (limiting) *Shape* of an ESD,  $\rho_{MP}(\lambda)$ , when the layer weight matrix has elements that are i.i.d. random from the Gaussian Universality class. In particular, the ESD will be MP when the matrix elements are drawn from a Normal distribution  $W_{i,j} \in N(0, \sigma^2)$ , e.g., as is typical at initialization, before NN training begins. Figure 2 (from Figure 4 of [?]) displays the MP distribution for different aspect ratios  $Q = \frac{N}{M}$  and variances  $\sigma^2$ . Notice that the *Shape* is characterized by a well defined, compact envelope with sharp edges.



(a) MP, varying  $Q = \frac{N}{M}$  (b) Tracy Widom fluctuations

Figure 2: MP distributions for different aspect ratios  $Q$  and variance scales  $\sigma^2$ , and an example of the finite-sized TW fluctuation  $\Delta_{TW}$ .

The MP distribution also predicts the *Scale* of an ESD, again when the layer weight matrix has elements that are i.i.d. random from the Gaussian Universality class. In particular, an MP distribution,  $\rho_{MP}(\lambda)$ , has very crisp, well-defined lower and upper bounds  $\lambda^-, \lambda^+$  [?], and (importantly) the upper bound  $\lambda^+$  exhibits finite-size Tracy-Widom (TW) fluctuations,  $\Delta_{TW}(\lambda)$ , which are on the order of  $\mathcal{O}(M^{-2/3})$ . Thus, any layer eigenvalue with *Scale* greater than this, i.e.,  $\lambda > [\lambda^+ + \Delta_{TW}(\lambda)]$ , is an “outlier” or a “spike.”

According to the HTSR phenomenology, these spikes carry significant generalizing information. (This is well-known for Bulk-Plus-Spike models [?], but the HTSR phenomenology generalizes this concept.) Relatedly, for layer matrices  $\mathbf{W}$  with aspect ratio  $Q > 1$  (i.e., rectangular matrices, where  $N > M$ ), MP RMT predicts there should be no zero eigenvalues, i.e.,  $\lambda_i > 0$ , for all  $i$ . Generally speaking, for well trained NNs, for layers with  $Q > 1$ , all eigenvalues are strictly larger than zero, i.e., well-trained layer weight matrices, with  $Q > 1$ , should have full rank and exhibit no “rank collapse.” HTSR places random Gaussian and “Bulk-Plus-Spike” matrices into the first two rows of Table 1. The essential feature of Gaussian random matrices is that their entries have no correlations. When some correlations are injected, a few large spike eigenvalues form, without otherwise disturbing the shape of the ESD. To really understand how individual NN layers converge, we need to understand when and why their ESDs become HT.

### 2.2.2 Heavy-Tailed Random Matrix Theory (HTRMT) and Power Law (PL) fits

For very well-trained NN layers, ESDs are *not* MP at all. Frequently, if not always, their ESDs are HT—and they are HT *because* they are strongly-correlated matrices. Importantly, they are *not* HT element-wise. Instead, their entries have a scale, and they have ESDs that are HT due to correlations learned during training. Existing theoretical approaches, including SLT and even StatMech, cannot readily model such strongly-correlated systems.<sup>10</sup>

Such strongly-correlated systems, however, do frequently arise in other, related scientific domains, including in the StatMech of self-organizing systems [?, ?], in electronic structure theory [?, ?, ?], and in quantitative finance [?, ?, ?]. In these (and other) domains, correlated systems frequently exhibit characteristic PL signatures; and it is common practice to *model* correlated systems as random (uncorrelated) systems by using HT statistics (e.g., Levy distributions or PL random matrices), fully understanding that such systems are by no means actually i.i.d. random. The HTSR phenomenology builds upon this longstanding practice by delimiting families of HT NN weight matrices based on the corresponding Universality classes of Pareto matrices.

<sup>10</sup>For example, such theoretical approaches typically deal better with *Scale* information (such as  $\lambda_{max}$ ) than with *Shape* information (such as  $\alpha$ ), e.g., by characterizing an “eigen-gap” separating large eigenvalues from “noise” [?] according to a noise plus low-rank perturbation model [?].



We explain briefly how to interpret Table 1 with respect to HTRMT. The 5+1 Phases of Training can be identified by fitting ESDs to MP or PL distributions, whichever gives the best fit, as shown in the last column. In case the PL distribution is a better fit, HTSR phenomenology treats the layer weight matrix as equivalent to an i.i.d. random matrix  $\mathbf{W}(\mu)$ , whose elements have been drawn from a Pareto distribution with exponent  $\mu$ .

**Heavy-Tailed Universality Classes of Random Pareto Matrices** For such an element-wise HT matrix, the theoretical *limiting* ESD of a Pareto matrix is also PL, which allows us to related the fitted PL  $\alpha$  with exponent  $\alpha = a\mu + b$ , to the Pareto exponent  $\mu$ . Ideally, for an infinite width matrix,  $a = \frac{1}{2}$  and  $b = 1$ , but due to finite-size effects, however, we have found we must take  $a \geq \frac{1}{2}$  and  $b \geq 1$ , giving

$$W_{i,j}(\mu) \sim \frac{C}{x^{\mu/2+1}}, \quad \rho(\lambda) \sim \lambda^{-(a\mu+b)}. \quad (6)$$

According to the above relation, we can use either the fitted PL exponent  $\alpha$ , or the Pareto exponent  $\mu$ , to index the HT Universality classes. Note, however, that the finite-size effects strongly depend on the aspect ratio  $Q = N/M$ , at least when applied to i.i.d random Pareto matrices, and the (Clauset MLE) PL fit may overestimate the  $\alpha$  of the ESD. Table 1 delimits the HT matrices into sub-categories (as shown in the bottom four rows) based on the behaviors of  $\alpha$  as a function of  $\mu$ .

Figure 3 illustrates how the fitted PL exponent  $\alpha$  corresponds to the actual Pareto exponent  $\mu$  for different aspect ratios  $Q = M/N$ . Figure 3(a) displays the ESDs of three different i.i.d.  $1000 \times 1000$  HT random matrices, with  $\mu = 1, 3, 5$ , on a Log-Log scale. Notice that smaller  $\mu$ , and therefore smaller  $\alpha$ , corresponds to heavier (i.e., larger) tails. Figure 3(b) shows how the empirically fit PL exponent  $\alpha$  can vary with the theoretical  $\mu$  for an associated  $\mathbf{W}(\mu)$ . For  $\mu < 2$  and  $Q = 1$ , the fitted  $\alpha$  follows the linear relation  $\alpha = \frac{1}{2}\mu + 1$ , albeit with some error. In contrast, for the more relevant  $\mu \in (2, 4)$  regime, the relation now depends far more strongly on the aspect ratio  $Q$ , and  $\alpha \in [2, 6]$ . For  $\mu > 4$ , the fitted  $\alpha$  saturates for each specific value of  $Q$ .

We emphasize that we only model the ESDs of the NN layer weight matrices using the same Universality class to that an associated with the ESD of an random, i.i.d, HT Pareto matrix. In fact, the elements  $W_{i,j}$  do not at all appear as if they have been drawn from a HT Pareto distribution, and, in contrast, are almost always well fit to a Laplacian distribution. Also, despite these strong finite-size effects, empirically one finds that the ESDs arising large, well trained, modern NNs can frequently be well fit to a PL (or TPL), and that the fitted  $\alpha \in [2, 6]$  for 80 – 90% of NN layers. Notably, we rarely find  $\alpha < 2$  in the best performing, open source, pretrained DNNs.

As there is *no* ground truth whatsoever as to the limiting spectral density of a strongly correlated NN weight matrix (especially without HT elements) the HTSR phenomenology uses Pareto matrices as a guide. However, as we will see in Section 3.3, this analogy should be treated with caution because there are cases where it breaks down.

No matter why matrix ESD is HT, it can be difficult to reliably estimating the  $\alpha$  parameter when the true  $\alpha$  is large. For Pareto matrices of the size investigated here, an observed  $\alpha$  above 6 is uninformative — the tail will decay very rapidly indeed, leaving very little of it to study. In this sense, the HT Universality classes are *larger* than the set of only strongly-correlated matrices or Pareto random matrices.

There is a particularly important boundary between Universality classes where  $\alpha = 2$ . Recall that one of the properties of power law distributions  $\rho(\lambda) \sim \lambda^{-\alpha}$  is that if  $\alpha < 2$ , then the variance of  $\rho(\lambda)$  is infinite. In such cases, the variance cannot be estimated empirically, making  $\rho(\lambda)$  in some sense *atypical*. This implies that the NN will have substantially greater difficulty in applying any further load to such a weight matrix. Thus, the value of  $\alpha = 2$  is a *critical value*. (See Figure 26 in Section 6.5 for an empirical study of this effect in a small MLP.)

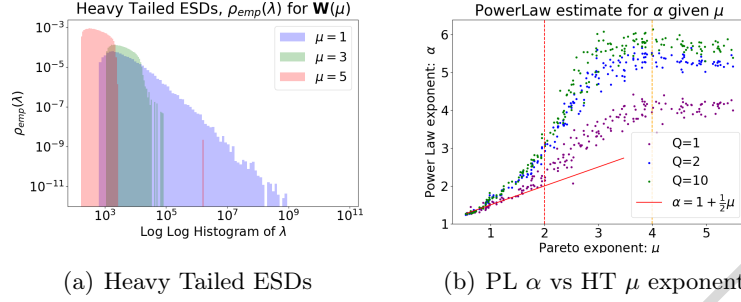


Figure 3: Comparison of ESDs and Power Law (PL) exponents  $\alpha$  from Heavy-Tailed (Pareto) weight matrices  $\mathbf{W}(\mu)$ . Subfigure (a) depicts 3 typical ESDs with Pareto exponent  $\mu = 1, 3, 5$ , each decreasing in *Shape* and *Scale*. Subfigure (b) shows how the exponent  $\alpha$  of the PL fit varies with  $\mu$ , with significant finite-size effects emerging for  $\mu > 2$  and  $\alpha > 2$ .

Smaller PL exponent  $\alpha$  values correspond to heavier tails,  $\rho_{tail}(\lambda)$ ; and the HTSR phenomenology observes that smaller PL exponents  $\alpha$  (at least for  $\alpha \in (2, 6)$ ) tend to correspond to better models. This is the key idea of the HTSR: the generalizing components of a layer matrix  $\mathbf{W}$  concentrate in larger singular vectors associated with the tail, and so that better models have more slowly-decaying (i.e., larger) ESD tails. This differs significantly than simply taking a general low-rank approximation to  $\mathbf{W}$ , where the rank is chosen without insight from the HTSR phenomenology. The SETOL theory formalizes this observation as a key assumption. We will revisit these model selection questions in Section 3.1 below.

### 2.3 Data-Free *Shape* and *Scale* Quality Metrics

The HTSR phenomenology provides quality metrics for both individual layers and (by averaging layers) for an entire NN model.

**Layer-wise Quality Metrics.** Using the HTSR phenomenology, we can define several other *Shape* and/or *Scale* based layer (quality) metrics. These are available in the `WeightWatcher` tool, and they work very well in practice.

- **Alpha ( $\alpha$ ):**  $\rho_{tail}(\lambda) \sim \lambda^{-\alpha}$ . A *Shape*-based quality metric.
- **LogSpectralNorm:**  $\log_{10} \lambda_{max}$ . A *Scale*-based quality metric.
- **AlphaHat ( $\hat{\alpha}$ ):**  $\alpha \log_{10} \lambda_{max}$ . A *Scale*-adjusted *Shape*-based quality metric.
- **Rand-Distance:**  $JSD[\rho^{emp} | (\rho_{rand}^{emp})]$ . A *Shape*-based, non-parametric quality metric, suitable for highly-accurate, epoch-by-epoch analysis.<sup>11</sup>
- **PL KS:**  $D_{KS}$ . The KS-distance, or quality-of-fit, of the PL fits. For transformers, foundation models, and large, complex, modern NNs, this is frequently an even better model quality metric than the  $\alpha$  of the PL fit itself.

<sup>11</sup>JSD is the Jensen-Shannon Divergence between the original ESD and the ESD of the layer weight matrix, randomized elementwise.



- **MP SoftRank:**  $\mathcal{R}_{MP}$ . The MP-SoftRank, defined in [?], can be used to identify problems such as when there is significant label or data noise that causes spuriously small  $\alpha$ , and also when it is difficult to fit a PL law.<sup>12</sup>

Each of these quality metrics provides a simple characterization of the *Shape* and/or *Scale* of the tail of the ESD of a given layer  $\mathbf{W}$ . These metrics are related to each other, and they have various trade-offs in practice [?, ?, ?]. Of particular interest here in our development of SETOL are the PL-based **WeightWatcher Alpha** and **AlphaHat** metrics.

**From Layer-wise Quality Metrics to Layer-Averaged Model Quality Metrics.** One can use the HTSR phenomenology to go beyond individual Layer Quality metrics, to construct model quality metrics by averaging Layer Quality metrics (over all layers that are not very small). Existing HTSR model quality metrics assume that all layers are statistically independent, so that the average model quality is just the average of the contributions from each weight matrix  $\mathbf{W}$ .<sup>13</sup> Given a Layer Quality metric,  $\bar{\mathcal{Q}}_L^{NN}(\mathbf{W})$ , one can define the *Model Quality*  $\bar{\mathcal{Q}}^{NN}$  metric for an entire model as

$$\bar{\mathcal{Q}}^{NN} := \prod_L \bar{\mathcal{Q}}_L^{NN}(\mathbf{W}), \quad (7)$$

a product of each independent Layer Quality  $\bar{\mathcal{Q}}_L^{NN}$ , and then consider the layer average as the log *Layer Quality*,

$$\log \bar{\mathcal{Q}}^{NN} = \frac{1}{N_L} \sum_L \log \bar{\mathcal{Q}}_L^{NN} = \langle \log \bar{\mathcal{Q}}_L^{NN} \rangle_{\bar{L}} \quad (8)$$

where  $\langle \dots \rangle_{\bar{L}}$  denotes the layer average.

In particular, prior work has used the following metrics:

- The layer-averaged model quality metric **Alpha**,  $\log \bar{\mathcal{Q}}^{NN} = \langle \alpha \rangle_{\bar{L}}$ , describes the *Shape* of the ESDs. One can use the averaged **Alpha** when studying a single model, and only varying the regularization hyperparameters, although **Alpha** also works very well as a model quality metric when comparing different transformer models [?].
- The layer-averaged model quality metric **LogSpectralNorm**,  $\log \bar{\mathcal{Q}}^{NN} = \langle \log \lambda_{max} \rangle_{\bar{L}}$ , describes the *Scale* of the ESDs. The averaged **LogSpectralNorm** does work as a model quality metric but not as well as **Alpha** (or **AlphaHat**). Notably, SLT predicts that smaller, not larger, **LogSpectralNorm** should be correlated with model quality; the opposite is observed in practice! This is because smaller layer  $\alpha$  generally, but not always, corresponds to larger  $\lambda_{max}$ .<sup>14</sup>
- The layer-averaged model quality metric **AlphaHat**,  $\log \bar{\mathcal{Q}}^{NN} = \langle \alpha \log \lambda_{max} \rangle_{\bar{L}} = \langle \hat{\alpha} \rangle_{\bar{L}}$ , incorporates both *Shape* and *Scale* information. This can compensate for anomalies that can arise when (say) comparing models of different sizes or model qualities [?] or when other issues cause unusually large  $\lambda_{max}$ . See Section 3.3.1).

<sup>12</sup>The **WeightWatcher** tool also implements the WW-SoftRank, which is like the MP-SoftRank, but replaces  $\lambda_{bulk}^+$  with  $\lambda_{rand}^{max}$ ; these are mostly equivalent for large matrices, but they can be different for very small matrices.

<sup>13</sup>This independence assumption, clearly a mathematical convenience, gets us closer to a workable theory. One could go beyond a “single layer theory” by adding in intra-layer correlations empirically. The **WeightWatcher** tool does support this, but doing so is outside the scope of this work.

<sup>14</sup>The **LogSpectralNorm** can exhibit a Simpson’s paradox when segmenting models by quality) [?]. Nevertheless, this metric may be useful when a PL fit can not be obtained, say, when  $N \gg M$  and  $M$  is very small, as with LSTMs, U-Net architectures, etc.

563 The layer-averaged **AlphaHat** model quality metric has been applied in a large meta-analysis  
564 of hundreds of SOTA pre-trained publicly-available NN models in CV and NLP [?, ?, ?, ?].  
565 Generally speaking, **HTSR** shape-based metrics, when used appropriately, outperform all other  
566 metrics studied (including those from **SLT**, and with access to the training/testing data,) for  
567 predicting the quality of SOTA pre-trained publicly-available NN models. The **HTSR** theory predicts  
568 that the best-performing NN models have layers with **Alpha**  $\in [2, 6]$ , and with  $\alpha = 2$  indicating  
569 optimal performance. Moreover, prior empirical results show that the **Alpha** and **AlphaHat** metrics  
570 can predict trends in the Quality (i.e., the Generalization Accuracy), of SOTA NN models—even  
571 *without access to any training or testing data* [?].

### 3 A Semi-Empirical Theory of (Deep) Learning (SETOL)

Based on prior empirical results, and the success of the **Alpha** and **AlphaHat** metrics that are based on the HTSR phenomenology, this leads to the deeper question:

*Why do the **Alpha** and **AlphaHat** metrics work so well as NN model quality metrics for SOTA NN models?*

That is, why do NN models with heavier-tailed layer ESDs tend to generalize better when compared to related models? Relatedly, can we derive these metrics from first principles? (If so, then under what conditions do they hold, and under what conditions do they fail?)

To answer these questions, we will derive a general expression for the Layer Quality,  $\bar{Q}$ , of a NN. Although many modern NNs have many layers, we adopt a single-layer viewpoint (like a matrix-generalized Student–Teacher) because in **SMOG** theory the multi-layer generalization can be factorized or approximated. For this, we will obtain by simple averaging our model quality metrics, under effectively a single layer approximation, that correspond to **Alpha** and **AlphaHat**.

In deriving these quantities, we will introduce to NN theory a new Semi-Empirical approach that combines techniques from **StatMech** and **RMT** in a novel way. The Layer Quality  $\bar{Q}$  will estimate the contribution that an individual NN layer makes to the overall quality of a trained NN model. In deriving  $\bar{Q}$ , we have discovered a new Layer Quality metric, called the **TRACE-LOG** condition, which indicates the generalizing components of the layer concentrate into a low-rank subspace (the *Effective Correlation Space*, or *ECS*). (Importantly, we have conducted detailed experiments to show that the key assumptions of our **SETOL** theory are valid (see Sections 6.2 and 6.3), and that the empirical estimates of the **SETOL** **TRACE-LOG** condition align remarkably well with predictions from the **HTSR** theory under Ideal conditions (see Sections 6.1). We also examine how the **HTSR** predictions (i.e., the HT PL exponent  $\alpha$ ) behave under non-Ideal conditions (see Sections 6.4 and 6.5).) In the following, we will outline key conceptual aspects of **SETOL**. In Section 3.1, ; In Section 3.2, ; and In Section 3.3.

#### 3.1 SETOL Overview

Our **SETOL** formulates a parametric expression for the Layer Quality  $\bar{Q}$  using a matrix-generalization of the classic Student–Teacher (ST) model from the Statistical Mechanics of Generalization (**SMOG**) theory of the 1990s [?, ?], evaluated using recent advances in the evaluation of so-called HCIZ random matrix integrals [?, ?, ?], such that the final expression for  $\bar{Q}$  can be written in terms of empirically measured statistical properties of the layer ESD. We summarize our basic approach here; see Section 5 for a detailed derivation, and see Section 6 for a detailed empirical analysis.

Following the Student–Teacher (ST) model [?], we first formulate the Generalization Error ( $\bar{\mathcal{E}}_{gen}^{ST}$ ) of the linear Perceptron (in the *Annealed Approximation*, and in the *High-Temperature* limit; see Section 4), and we then generalize this to the case of a NN ( $\bar{\mathcal{E}}_{gen}^{ST} \rightarrow \bar{\mathcal{E}}_{gen}^{NN}$ ), so that we can analyze the Quality of each layer. For the Perceptron, the Generalization Error is an Energy, given as  $\bar{\mathcal{E}}_{gen}^{ST} := \langle 1 - R \rangle_s^\beta$ , where  $R$  is the ST vector overlap, and  $\langle \dots \rangle_s^\beta$  is a *Thermal Average* (defined in Section 4.2), a Boltzmann-weighted average. In this case, the model Quality,  $\bar{Q}^{ST}$  is exactly the AA, high-T Average Generalization Accuracy  $\bar{Q}^{ST} := 1 - \bar{\mathcal{E}}_{gen}^{ST} = \langle R \rangle_s^\beta$ . For an MLP or general NN, each layer Energy is associated with a Layer Quality  $\bar{Q}$ , which we identify as the average contribution an individual layer makes to the overall generalized accuracy. (i.e  $1 - \bar{\mathcal{E}}_{gen}^{NN}$ ) for a multilayer perceptron (MLP).

For technical reasons (below), we will seek the *Layer Quality (Squared)*  $\bar{Q}^2$ , which is defined as

the Thermal Average of the matrix-generalized overlap ( $\text{Tr}[\mathbf{R}^2]$ ),

$$\bar{Q}^2 := \langle \text{Tr}[\mathbf{R}^2] \rangle_{\mathbf{S}}^{\beta} \quad (9)$$

where  $\mathbf{R}^2$  can be thought of as a Hamiltonian for the Quality-Squared ( $\mathbf{H}_{\bar{Q}^2} = \mathbf{R}^{\top} \mathbf{R}$ ).

In Eqn. 9, the so-called *Teacher* ( $T$ ) is the NN model under consideration, and  $\mathbf{R} := \frac{1}{N} \mathbf{S}^{\top} \mathbf{T}$  denotes the ST overlap operator between the Teacher layer weight matrix  $\mathbf{T}$  and a similar *Student* ( $S$ ) layer weight matrix  $\mathbf{S}$ . The notation  $\langle \dots \rangle_{\mathbf{S}}^{\beta}$  denotes a Thermal Average over all Student weight matrices  $\mathbf{S}$  that resemble the Teacher weight matrix  $\mathbf{T}$ . By “resemble”, the SETOL approach assumes that the ESD of  $\mathbf{S}$  has the same *limiting* form as  $\mathbf{T}$ , placing them in the same HTSR Universality class. This is made more precise below.

Let us now express the average matrix-matrix overlap  $\mathbf{R}$  in squared form using:

$$\begin{aligned} \text{Tr}[\mathbf{R}^2] &:= \text{Tr}[\mathbf{R}^{\top} \mathbf{R}] \\ &= \frac{1}{N^2} \text{Tr}[\mathbf{T}^{\top} \mathbf{S} \mathbf{S}^{\top} \mathbf{T}] = \frac{1}{N} \text{Tr}[\mathbf{T}^{\top} \mathbf{A}_2 \mathbf{T}] \end{aligned} \quad (10)$$

where  $\mathbf{A}_2$  is the  $N \times N$  form of the Student correlation matrix,  $\mathbf{A}_2 := \frac{1}{N} \mathbf{S} \mathbf{S}^{\top}$ . We will also define the  $M \times M$  matrix,  $\mathbf{A}_1 = \frac{1}{N} \mathbf{S} \mathbf{S}^{\top}$ , used later.

As explained in Section 4.2, this Quality-Squared is more readily obtained as the derivative of the Layer Quality-Squared Generating Function,  $\beta \Gamma_{\bar{Q}^2}^{IZ}$ , defined as

$$\bar{Q}^2 := \frac{1}{\beta} \frac{\partial}{\partial N} \lim_{N \gg 1} \beta \Gamma_{\bar{Q}^2}^{IZ} \quad (11)$$

$\beta \Gamma_{\bar{Q}^2}^{IZ}$  is essentially ( $\beta$  times) a *Free Energy* for the (approximate) Layer Quality-Squared. (see Section 4, and the Appendix, Section A.3). For more details, see Section 5, and the Appendix, Section A.2).

We can write  $\beta \Gamma_{\bar{Q}^2}^{IZ}$  as an HCIZ Integral,

$$\beta \Gamma_{\bar{Q}^2}^{IZ} = \ln \int d\mu(\mathbf{S}) \exp(N\beta \text{Tr}[\mathbf{T}^{\top} \mathbf{A}_2 \mathbf{T}]), \quad (12)$$

The SETOL approach then seeks to express  $\beta \Gamma_{\bar{Q}^2}^{IZ}$  in Eqn. 12 as an HCIZ integral (and in the large- $N$  limit) [?, ?, ?]. We evaluate this at large- $N$ , and write

$$\beta \Gamma_{\bar{Q}^2, N \gg 1}^{IZ} := \lim_{N \gg 1} \beta \Gamma_{\bar{Q}^2}^{IZ} \quad (13)$$

With these definitions in place, moving forward, the following key assumptions, which can be tested empirically, must hold:

- **The Effective Correlation Space (ECS) Condition.** The generalizing components of the Student (and Teacher) layer weight matrices concentrate into a lower rank subspace—the ECS—spanned by the eigenvectors associated with the (heavy) tail of the layer ESD  $\rho_{\text{tail}}(\lambda)$ , such that the test error can be reproduced with only these components. We write  $\tilde{\mathbf{A}}$  to denote the projection of the correlation matrix  $\tilde{\mathbf{A}} := \mathbf{P}_{\text{ecs}} \mathbf{A}$ , onto this subspace, now with rank  $\tilde{M} \ll M$ . This restricts the measure  $d\mu(\mathbf{A})$  to the ECS, ( $d\mu(\mathbf{A}) \rightarrow d\mu(\tilde{\mathbf{A}})$ ). This assumption will empirically be examined using real-world Teacher weight matrices  $\mathbf{T} = \mathbf{W}$  in Section 6.2.

644 • **The TRACE-LOG Condition.** The Student correlation matrix  $\tilde{\mathbf{A}}_1$  (when properly normalized)  
 645 satisfies the condition that  $\text{Tr}[\ln \tilde{\mathbf{A}}_1] = \ln \det(\tilde{\mathbf{A}}_1) = 0$ , so that the change of measure  
 646  $d\mu(\mathbf{S}) \rightarrow d\mu(\tilde{\mathbf{A}})$  is Volume Preserving. This condition is derived explicitly in terms of  $\tilde{\mathbf{A}}_1$ ,  
 647 and therefore will hold for  $\tilde{\mathbf{A}}_2$  (and the Teacher Correlation matrix,  $\tilde{\mathbf{X}} = \frac{1}{N} \mathbf{T}^\top \mathbf{T}$ ). Practically,  
 648 this implies that the  $\tilde{M}$  eigenvalues  $\tilde{\lambda}$  of the tail of the ESD must satisfy  $\sum_{i=1}^{\tilde{M}} \ln \tilde{\lambda}_i \approx 0$ .  
 649 Experiments will test this assumption explicitly in Section 6.3.

650 Remarkably, both conditions hold best empirically when the HTSR PL quality metric  $\alpha \gtrsim 2$  is Ideal.  
 651 Motivated from these empirical observations, we have:

652 •  $\beta \Gamma_{\tilde{\mathcal{Q}}^2, N \gg 1}^{IZ}$  is expressed as an HCIZ integral, at large- $N$ . We have

$$\beta \Gamma_{\tilde{\mathcal{Q}}^2, N \gg 1}^{IZ} = \lim_{N \gg 1} \ln \int d\mu(\tilde{\mathbf{A}}) \exp(\beta \text{Tr}[\mathbf{T}^\top \tilde{\mathbf{A}}_2 \mathbf{T}]) \quad (14)$$

653 where measure  $d\mu(\tilde{\mathbf{A}})$  lets us average over all Student Correlation matrices  $\tilde{\mathbf{A}}_2$  which lie in  
 654 the ECS space and which “resemble” the Teacher, where by “resemble” we mean that they  
 655 share the same form of the tail of their limiting ESDs, i.e.,  $\rho_{\tilde{\mathbf{A}}}^\infty(\lambda) \sim \rho_{\tilde{\mathbf{X}}}^\infty(\lambda)$ .

656 • **The Layer Quality (Squared)  $\tilde{\mathcal{Q}}^2$  is a Norm Generating Function.** The final  
 657 expression for  $\tilde{\mathcal{Q}}^2$  can be written as the derivative of  $\beta \Gamma_{\tilde{\mathcal{Q}}^2, N \gg 1}^{IZ}$  as

$$\tilde{\mathcal{Q}}^2 = \frac{1}{\beta} \frac{\partial}{\partial N} \beta \Gamma_{\tilde{\mathcal{Q}}^2, N \gg 1}^{IZ} = \sum_{i=1}^{\tilde{M}} \mathcal{G}(\lambda_i) \quad (15)$$

658 where  $\mathcal{G}(\lambda_i)$  is a *Norm Generating Function*, and is defined as the integrated *R-transform*  
 659  $R(z)$  of the Teacher layer ESD (where  $z \in \mathbb{C}$ ), such that  $\mathcal{G}(\lambda) := \int_{\lambda_{min}^{ECS}}^{\lambda} R(z) dz$  and  $\lambda_{min}^{ECS}$   
 660 encapsulates the ECS (and selects the desired branch-cut of  $R(z)$  so that it is both single-valued  
 661 and analytic).

662 To apply the theory, one must choose and R-transform  $R(z)$  for the Teacher that models the  
 663 tail of the ESD  $\rho_T^{emp}(\lambda)$ , and which can be parameterized by some measurable property. This may  
 664 include the number of Spikes  $\lambda^{spike}$ , the fitted PL exponent  $\alpha$ , the maximum eigenvalue  $\lambda_{max}$ , or  
 665 even the entire tail  $\rho_T^{tail}(\lambda)$ . This may be a formal expression, a computational procedure, or some  
 666 combination.

667 To integrate  $R(x)$ , however, to have a physically meaningful result, one must ensure that  $R(z)$   
 668 is both analytic and single-valued on domain of interest, namely, the ECS (and therefore the (PL)  
 669 tail of the ESD),  $z \geq \lambda_{min}^{ECS}$ . Because the ESD is frequently Heavy-Tailed (HT), this R-transform  
 670  $R(z)$  may have a branch-cut, and it is expected that this will occur at  $z \leq \lambda_{min}^{ECS}$ , corresponding  
 671 roughly at or before the start of the ECS. In a sense, selecting the branch-cut  $R(z)$  forces one to  
 672 define the ECS.

673 To complete the theory, we will also show that the HTSR PL Layer Quality metrics **Alpha** ( $\alpha$ )  
 674 and **AlphaHat** ( $\hat{\alpha}$ ) can be formally derived directly from the SETOL Layer Quality  $\tilde{\mathcal{Q}}$  by selecting  
 675 the appropriate R-transform  $R(z)$ . In Section 5.4 we provide several possible model  $R(z)$  and the  
 676 resulting Layer Quality  $\tilde{\mathcal{Q}}$ .

677 **Renormalization Group Effective Hamiltonian** The formulation of SETOL closely paral-  
 678 lels the construction of an Effective Hamiltonian  $\mathbf{H}_{\tilde{\mathcal{Q}}^2}^{ECS}$  via the Wilson Exact Renormalization  
 679 Group (RG) approach. Consider a *bare* Hamiltonian  $\mathbf{H}_{\tilde{\mathcal{Q}}^2}$  for the Layer Quality-Squared, defined as  
 680  $\mathbf{H}_{\tilde{\mathcal{Q}}^2} := \mathbf{R}^\top \mathbf{R}$ . We can express Eqn.12 in terms of this bare Hamiltonian  $\mathbf{H}_{\tilde{\mathcal{Q}}^2}$ , and rewrite Eqn.14

in terms of an *renormalized* Effective Hamiltonian  $\mathbf{H}_{\bar{Q}^2}^{ECS}$  that spans the Effective Correlation Space (ECS). Formally, we have:

$$\ln \int d\mu(\mathbf{S}) e^{N\beta \text{Tr}[\mathbf{H}_{\bar{Q}^2}]} \xrightarrow{RG} \lim_{N \gg 1} \ln \int d\mu(\tilde{\mathbf{A}}) e^{N\beta \text{Tr}[\mathbf{H}_{\bar{Q}^2}^{ECS}]} \quad (16)$$

where the RG transformation is defined by the Scale-Invariant TRACE-LOG condition, applied at large- $N$ , and where  $\mathbf{H}_{\bar{Q}^2}^{ECS}$  is defined implicitly through result for  $\bar{Q}^2$  (Eqn. 15). The result is, formally, a sum of the integrated R-transforms  $\mathcal{G}(\lambda_i)$ .<sup>15</sup> [For the conclusions:] The presence of the branch-cut in  $R(z)$  suggests a situation that is similar in spirit to the RG solution characterizing a phase transition, here with a critical exponent of  $\alpha = 2$ , and the phase boundary being between the Heavy-Tailed (HT) and the Very Heavy-Tailed (VHT) phase of learning of the HTSR theory. This observation strengthens our argument that the HTSR HT and HVT phases indeed are analogous to the generalizing and overfit phases, resp., of the classical SMOG theories of NN learning.

### 3.2 Comparing SETOL with HTSR: Conditions for Ideal Learning

The SETOL approach establishes a starting point for developing a first-principles theory for modern NNs. Among other things, by connecting with the HTSR phenomenology, it lets us identify conditions for an Ideal state of learning for an individual NN layer, under the Single Layer Approximation. By Ideal, we mean that the layer is being used most effectively i.e., in some sense it is at its optimal data load, and thus it is conjectured to result in the best model quality.

**The Ideal State of Learning** is conjectured to be characterized by the following three conditions:

1. the tail of ESD,  $\rho_{tail}^{emp}(\lambda)$ , can be well fit to a PL of  $\alpha \approx 2$ :  $\rho_{tail}^{emp}(\lambda) \sim \lambda^{-2}$ ;
2. the eigenvalues in the tail,  $\lambda_i$ , satisfy the TRACE-LOG condition:  $\sum_i \ln \lambda_i = 0$ ; and
3. the generalizing components of the layer concentrate in the singular vectors associated with the tail of the ESD, (Effective Correlation Space).

In Section 6, we will test and justify this conjecture.

These claims are fundamentally about NN learning itself. They are motivated by our formulation of the SETOL approach in our search for a practical predictive theory behind the HTSR phenomenology. When (1) and (2) conditions hold for any layer, we conjecture that (3) holds as well. Moreover, when (1–3) hold for all layers, we conjecture the NN has the lowest Generalization Error (and highest Model Quality) possible for given model architecture and dataset.

Previous results have shown that the HTSR quality metrics (Alpha, AlphaHat, etc.) correlate very well with reported test accuracies, as well as model quality on epoch-by-epoch basis. These results hold because, as indicated by the HTSR theory, the PL exponent  $\alpha$  characterizes both the quality of the layer and provides an after-the-fact measure of the amount of regularization.<sup>16</sup> However, the HTSR approach says nothing about the SETOL TRACE-LOG condition; and neither does the SETOL approach require a minimum of  $\alpha = 2$  to obtain the best model quality, as observed by the HTSR phenomenology. Remarkably, we can show that (1) and (2) do hold *simultaneously*, both

<sup>15</sup>In a sense, this result resembles (a non-perturbative form of) the Linked Cluster Theorem in that the log Partition Function is expressed as a sum of integrated matrix-generalized cumulants.

<sup>16</sup>By “after-the-fact”, we mean that it provides a measure of the regularization in a layer, along the lines of the self-regularization interpretation of HTSR Theory [?]. However, we do *not* recommend that it be used as an explicit regularization parameter. Informally, this is since the “easiest way to obtain HT ESDs is to make weight matrices HT element-wise; but this is *not* what is observed in practice, and thus this is precisely *not* what HTSR Theory and our new SETOL approach are designed to model.



in carefully designed experiments on a small model, as well as for many pre-trained, high quality open-source models (such as VGG, ResNet, Llama, Falcon, etc).

HTSR theory, however, has been developed as a phenomenology describing the best-trained, most accurate open-source models available. As such, it may be biased towards such models, and it may not describe less optimal learning scenarios. The keys goals of this work are to derive independent conditions, both theoretical and experimental, that can identify the conditions for Ideal Learning, and to stress-test these conditions in carefully designed, reproducible experiments.

### 3.3 Detecting Non-Ideal Learning Conditions

The HTSR phenomenology posits that SGD training reduces the *Training Error* by accumulating correlations into the large eigenvalues in NN layer weight matrices  $\mathbf{W}$  such that they *self-organize* into a HT with a PL signature, and that this successful self-organization leads to good model quality. Conversely, it also posits that when training has gone awry in some way, the resulting ESD,  $\rho^{emp}(\lambda)$ , will be deformed in some way. In many practical situations, there can be other, competing factors that give rise to large eigenvalues ( $\lambda > \lambda^+$ ) that do not contribute to the generalization capacity of the model, and, consequently, can affect the Scale (i.e., the largest eigenvalue(s)  $\lambda_{max}$ ) and Shape (i.e., the PL exponent  $\alpha$ , or goodness of fit  $D_{KS}$ ) of the layer ESDs. These large  $\lambda$  could be Dragon Kings, *Correlation Traps*, or some other anomaly.

In order to apply the HTSR phenomenology most effectively, one must be able to identify various spurious factors and distinguish real correlations from any other large eigenvalues, including the effects of both extreme eigenvalues  $\lambda$ , individual matrix elements  $W_{i,j}$ , and rank-1 perturbations in  $\mathbf{W}$ . In one case the ESD is HT primarily due to correlations that help the model generalize, whereas in another when the ESD may be more HT than expected due to suboptimal training, mis-labeled data, etc. In extreme cases, spurious eigenvalues can push the weight matrix into the Very Heavy-Tailed Universality class (i.e.,  $\alpha < 2$ . See Table 1, Section 2), or disrupt the formation of a HT, resulting in a poor PL fit, undermining the core proposition of the HTSR approach.

When training a model with SGD, one may only achieve a sub-optimal result when using overly large learning rates / small batch sizes, (see Section 6), from poor hyper-parameter settings, or simply because direct regularization fails. In such cases, the HTSR approach allows one to detect potential problems by looking for large eigenvalues *not* resulting from correlations.[?]

Importantly, in the context of the SETOL theory, we can now identify such empirical anomalies due to atypical layer weight matrices  $\mathbf{W}$ , a key factor when models break down. The SETOL approach formalizes the empirical HTSR phenomenology, but, in doing so, assumes that the underlying layer effective correlation matrix  $\tilde{\mathbf{X}}$ , is typical, meaning that it can describe out-of-sample / test data. Conversely, if the underlying weight matrix  $\mathbf{W}$  is atypical, then it is in some sense overfit to the training data and can not fully represent out-of-sample / test data. Consequently, when  $\mathbf{W}$  is atypical, we argue that we can observe this, either in the ESD  $\rho^{emp}(\lambda)$  directly (i.e., when  $\alpha < 2$ ), and/or having unusually large matrix elements  $W_{ij}$ .

We conjecture such sub-optimal results, and particularly those occurring from overfitting, actually arise when the underlying layer weight matrix  $\mathbf{W}$  is atypical in some way, in analogy to the results from the classic SMOG theory (see Section 4.1), and, importantly, that we can use the SETOL approach to detect when  $\mathbf{W}$  is atypical and therefore a layer is overfit in some way.

Here, we identify two specific cases of atypical weight matrices—Correlation Traps and *Over-Regularization*.<sup>17</sup>

**3.3.1 Correlation Traps.**  $\mathbf{W}$  is atypical in that  $\mathbf{W}$  exhibits an anomalously large mean ( $\bar{\mathbf{W}}$ ).

<sup>17</sup>Later, in Section 6, we will show that we can systematically induce both phenomena and observe their effects on the HTSR HT PL metric  $\alpha$  and the SETOL TRACE-LOG condition.

We can observe these by randomizing the layer weight matrix,  $\mathbf{W} \rightarrow \mathbf{W}^{rand}$ , and then looking for eigenvalues that extend significantly beyond the MP edge of the random bulk (i.e., Spikes).. We call such random spikes *Correlation Traps*, denoted as  $\lambda_{trap}$ , because they appear, in some extreme cases, to be associated with distorted ESDs and, importantly, lower test accuracies. Examples of Correlation Traps are shown in Section 6.4.

**3.3.2 Over-Regularization.**  $\mathbf{W}$  is atypical in that  $\mathbf{W}$  exhibits an anomalously large variance ( $\sigma^2(\mathbf{W})$ ). We can observe this when the layer  $\alpha < 2$ . Also, since  $\mathbf{Alpha}$  a measure of implicit regularization, we say the layer with  $\alpha < 2$  is *Over-Regularized*. In particular, when one layer is undertrained, having  $\alpha > 6$ , it appears that other layers may become overtrained to compensate, and this can be seen with having  $\alpha < 2$ . These effects are studied in Section 6. Additionally, we also observe that when evaluating the SETOL TRACE-LOG condition, when  $\alpha < 2$ , then  $\Delta\lambda_{min} < 0$  (see Section 6.3).

### 3.3.1 Correlation Traps

The first way we identify  $\mathbf{W}$  as atypical is when it has an anomalously large mean ( $\bar{\mathbf{W}}$ ); detecting this in general, however, requires more than just examining which elements  $W_{i,j}$  are anomalously large.

The HTSR phenomenology states that NNs generalize better when their layers ESDs are more HT—precisely because the tail eigenvalues arise from correlations in the weight matrices. So one way is to identify *atypicality* is to look for large eigenvalues that do not arise from correlations in  $\mathbf{X}$ , but, rather, from one or a few spuriously large matrix elements  $W_{i,j}$  and/or rank-1 perturbations in  $\mathbf{W}$ . We call these eigenvalues Correlation Traps, denoted by  $\lambda_{trap}$  (i.e., see Section 6.4).

Indeed, if we randomize  $\mathbf{W}$  elementwise, i.e  $\mathbf{W} \rightarrow \mathbf{W}^{rand}$ , we expect the  $W_{i,j}^{rand}$  matrix elements to be i.i.d and with a small mean (unless something odd happens during SGD training). Likewise, we expect the singular values of  $\mathbf{W}^{rand}$  to follow the MP distribution, to within finite-size / TW fluctuations. If we observe an eigenvalue  $\lambda_{trap}$  extending beyond the MP bulk region,  $\lambda_{trap} > \lambda_{bulk}^+$ , then the mean  $W_{i,j}$  matrix element will also be anomalously large, and we can identify  $\mathbf{W}$  as atypical. We must be careful, however, as we do not fully understand the origin of these atypicalities and do not claim that every one is associated with suboptimal generalization.

By a *Correlation Trap*, we mean that some anomaly in the training of  $\mathbf{W}$  resulted in one or more spuriously large eigenvalues  $\lambda_{trap}$  in  $\mathbf{W}^{rand}$ , and that whatever caused them also may, in some pronounced cases, tend to “trap the correlations in  $\mathbf{X}$  itself, preventing them from coalescing into a well defined PL Heavy Tail, or otherwise distorting the ESD. Whether they are a signature of training gone wrong, or whether they distort the dynamics of the tail correlations simply by being there, Correlation Traps can be expected to alter the shape of the ESD, reducing the quality of the PL fits, and sometimes producing spurious  $\alpha$  values.

Why would such anomalies occur in a NN? It is conceivable that SGD will, when it fails to find usable correlations, instead produce spurious correlations in the form of large elements and/or rank-1 perturbations. Also, the matrix itself may simply undergo an innocuous mean-shift because the mean is not explicitly controlled during training. Here, mean-recentering may be beneficial.<sup>18</sup>

We will see, below in Section 6, that we can induce a Correlation Trap both by shrinking the batch size, or, equivalently, increasing the learning rate, and that this is associated with degraded model performance and small  $\alpha$ . We seek to identify specific ways of identifying such traps because we reason that the self-organization of correlations may be disrupted by the presence of foreign large eigenvalues – or that failed learning may produce them as a by-product.

<sup>18</sup>Similarly, when training NNs, frequently the weight matrices [?] or activations [?] may need to be clipped during training to ensure good results.



**Detecting Correlation Traps with RMT.** RMT suggests that when a matrix  $\mathbf{W}$  has unusually large elements  $W_{ij}$ , then the ESD will have one or more large eigenvalues  $\lambda$  lying outside the bulk edge  $\lambda_{bulk}^+$  of the ESD, as predicted by MP theory. One can detect these so-called Correlation Traps in a weight matrix  $\mathbf{W}$  by performing the following:

1. randomize  $\mathbf{W}$  element-wise to obtain  $\mathbf{W}^{rand}$ ;
2. compute the ESD for  $\mathbf{W}^{rand}$ ; and
3. look for large eigenvalues  $\lambda_{trap} \gg \lambda_{bulk}^+$ .

**WeightWatcher** looks for Correlation Traps ( $\lambda_{trap}$ ) in the ESD of the randomized  $\mathbf{W}^{rand}$ , that are larger than  $\lambda_{trap} > \lambda_{bulk}^+ + \Delta_{TW}$ , where  $\lambda_{bulk}^+$  is the MP bulk edge  $\Delta_{TW}$  are the associated finite-size Tracy Widom (TW) fluctuations. This procedure detects *any* anomaly in the matrix weights that produce spuriously large eigenvalues. It is implemented in **WeightWatcher** (using the randomize option), which was used to generate the plots in Figure 4.

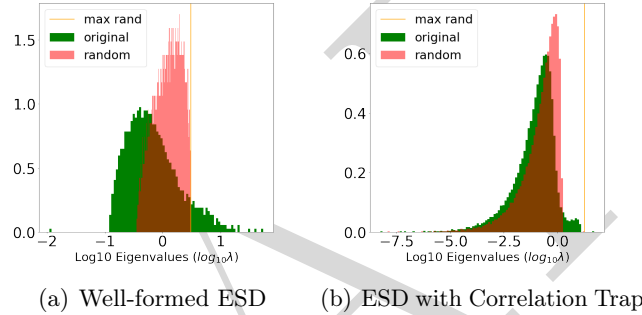


Figure 4: Comparison of a well-formed, Heavy-Tailed ESD (a) to one with a Correlation Trap (b), in the VGG16 model (FC2 layer)

See Figure 4(a), which displays the (log)-ESD of a typical SOTA NN layer  $\mathbf{X}$  (green), i.e., on a log-linear scale, along with the (log)-ESD of that layer after randomizing it element-wise (red). The two ESDs differ substantially: the ESD of the original weight matrix  $\mathbf{W}$  (green) is very HT, whereas the ESD of the randomized weight matrix  $\mathbf{W}^{rand}$  (red) is an MP (and as predicted by the MP RMT). The orange line corresponds to the maximum eigenvalue of the randomized ESD. Note that it is at the MP bulk edge of the red plot, indicating that this ESD is not affected by unusually large elements or other weight anomalies. Here, we say that the ESD of  $\mathbf{X}$  is HT, and that  $\mathbf{W}$  is not HT element-wise. HTSR says this layer is well-trained.

Contrast this with Figure 4(b), which displays the (log)-ESD of a NN layer with a Correlation Trap. The ESD of  $\mathbf{X}$  (green) is weakly HT, but it looks nothing like the ESD in Figure 4(a). In fact, it looks very much like the ESD of the randomized weight matrix  $\mathbf{W}^{rand}$  (red), except for a small shelf on the right. The orange line again corresponds to the maximum eigenvalue of the randomized ESD, and this is just past this shelf. Relative to the randomized ESD, this line depicts (an) unusually large element(s)—or, equivalently, a rank-1 perturbation of  $\mathbf{W}^{rand}$ . By a Correlation Trap, we mean that some anomaly in the elements of  $\mathbf{W}$  tends to “trap” the ESD of  $\mathbf{X}$ , concentrating the correlations in  $\mathbf{X}$  into the small shelf of density around the orange line. HTSR says this layer is not well-trained because it does not have a good PL fit.

### 3.3.2 Over-Regularization

The second way we identify  $\mathbf{W}$  as atypical is when it has an anomalously large variance ( $\sigma^2(\mathbf{W})$ ).

The SETOL theory – a single-layer theory of learning – casts the training of a NN layer in terms of how the correlations concentrate into the layer Effective Correlation Space (ECS), and becomes exact when the TRACE-LOG condition is satisfied. Analogously, the HTSR theory – a single layer phenomenology of learning – casts training of an N layer by fitting its ESD to a PL, and noting that the PL exponent  $\alpha$  measures the amount of implicit regularization in the layer. Comparing the two approaches, we see that smaller  $\alpha$  corresponds to the correlations concentrating into a low-rank ECS. In general, and likewise, the more the weight matrix correlations concentrate into a low-rank ECS, the better the layer has been regularized. A natural question arises then, namely, can a layer be *Over-Regularized* and can we detect this? and in large, Empirically, we do indeed observe that over the course of training,  $\alpha$  decreases, (See Figure 26 (a), Section 6.5,) and that the models predictions are concentrated into the ECS, (See Section 6.2). Thus, we also interpret  $\alpha$  and ECS concentration to be measures of learning itself, meaning that NNs are self-regularizing [?].

Importantly, however, the HTSR phenomenology indicates that Alpha usually lies in the Fat-Tailed Universality class, such that  $\alpha \in [2, 6]$ . When  $\alpha < 2$ , the ESD is Very Heavy Tailed (VHT), and, also, this indicates that  $\mathbf{W}$  has an anomalously large variance. That is,  $\mathbf{W}$  is atypical. Occasionally, but very infrequently, we do observe  $\alpha < 2$ , and in large, production quality models (like Llama). Interestingly, we also observe that, frequently, when the HTSR  $\alpha < 2$ , the SETOL TRACE-LOG condition holds fairly well. This is further explored in Section 6

We have applied the WeightWatcher tool to have examined dozens of modern, very large NNs; of particular interest are the so-called Large Language Models (LLMs) that have revolutionized the field of AI. To that end, in Fig. 5, we present the WeightWatcher layer Alpha metrics for the Falcon-40b and the Llama-65b LLMs.<sup>19</sup>

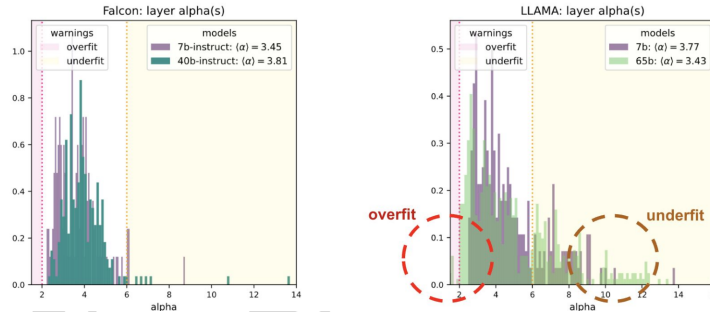


Figure 5: Falcon vs Llama

For the Falcon-40b model, all of the layer Alpha range between  $\alpha \in [2, 6]$ , and therefore lie in the Fat-Tailed Universality class (in Table 1) and are well-fit. In contrast, looking at the Llama-40b layer Alpha, very many have  $\alpha > 6$ , indicating these layers are under-fit, and while a few have  $\alpha < 2$ , suggesting these are over-fit. Finally, there are more layers with  $\alpha \sim 2$  in Llama-65 vs Falcon-40b.

The observations on Llama-2 suggest that the layers with  $\alpha \leq 2$  are compensating for the layers with  $\alpha > 6$ , and yielding suboptimal performance for the Llama-65b architecture. Based on these observations, we hypothesize that, in a multi-layer-perceptron (MLP), when one layer does not or can not learn, then other layers will have to compensate, and will be overloaded with the training data, leading to  $\alpha < 2$ , and even the TRACE-LOG condition  $\Delta\lambda_{min} < 0$ .

In Section 6.5, we will test this hypothesis. By reducing the trainable parameters in a small MLP, we can simulate the situation seen above in the Llama-65b model, and observe the formation of a Very Heavy Tailed (VHT) ESD in the dominating layer weight matrix. Overloading results

<sup>19</sup>Similar results are found for the larger, more modern Llama models, and can be found on the WeightWatcher website[?]

869 from having too few parameters for the complexity of the task. Adding more data increases the  
870 load up to the total complexity of the task itself.

871 Moreover, we will also argue that in our experiments, the model enters a kind of glassy meta-  
872 stable phase, similar to the kinds of phases predicted by classic **StatMech** theories of learning [?]  
873 (described below). Section 6.5 will explore how far we can push the analogy of glassy systems in  
874 our experiments to stress test the **SETOL** approach. In particular, we will see effects such as the  
875 slowing down of its dynamics, leading to a kind of hysteresis, specific to the under-parameterized  
876 regime.

DRAFT

## 4 Statistical Mechanics of Generalization (SMOG)

In this section, we review the **StatMech** approach to learning: both to understand how it is usually applied in Statistical Mechanics of Generalization (SMOG) theory; and to understand how our Semi-Empirical approach in SETOL is similar to and different from the traditional approach. We will also obtain an expression for the Generalization Accuracy (or Model Quality  $\bar{Q}^{ST}$ ) for the classic Student-Teacher (ST) model of the Linear Perceptron (in the AA, and at high-T), as described in [?, ?]. In Section 5, we will generalize this to a Layer Quality metric,  $\bar{Q}$ , for a layer in a general Multi-Layer Perceptron (MLP), i.e.,  $\bar{Q}^{ST} \rightarrow \bar{Q}$ , so that  $\bar{Q}$  can then be expressed in terms of the ESD of the NN layer.

**Outline.** Here is an outline of this section.

- **Approaches to the SMOG.** In Section 4.1, we explain the mapping from the **StatMech** theory of disordered systems to the **StatMech** theory of NN learning (SMOG); and how our Semi-Empirical approach (SETOL) is similar to and different from the traditional approach.
- **Mathematical Preliminaries.** In Section 4.2, we review mathematical details of **StatMech**, providing definitions and detailed derivations of quantities and expressions necessary later.
- **Student-Teacher Model.** In Section 4.3, we discuss the setup of the Student-Teacher (ST) model as a general means to estimate the Average Generalization Error empirically. First, in subsection 4.3.1, we describe the ST setup with an operational analogy. Then, in subsection 4.3.2, we derive the (new) result for the ST Model Quality,  $\bar{Q}^{ST}$ , using the setup of the classic (ST) model for the Generalization Error (and accuracy) of the Perceptron (in the AA, and at high-T).

Additional information can be found in the Appendix.

- **Symbols and Equations.** In Section A.1, we summarize the important symbols and key results, including the dimensions of the vectors and matrices, different notations for energies, and key equations.
- **Summary of the SMOG.** In Section A.2, we provide a more detailed analysis of the results we derive in Section 4.3.2. In particular, in Section A.2.1, we repeat the derivations of the ST Generalization Error  $\bar{\mathcal{E}}_{gen}^{ST}$  and related quantities (from Section 4.2), using more concrete notation to align with [?, ?]; and in Section A.2.2, we use this to derive the matrix-generalization of the ST Annealed Error Potential  $\epsilon(R)$  (as well as the normalization for the weight matrices, necessary for later).

### 4.1 StatMech: the SMOG approach and the SETOL approach

In this subsection, we review the basic **StatMech** setup necessary to understand SMOG theory as well as SETOL. This theory was developed in the 1980s and 1990s [?, ?, ?, ?, ?].

**Traditional SMOG theory.** In traditional SMOG theory, one maps the learning process of a NN to the states and energies of a physical system. The mapping is given in Table 2. SMOG theory models the SGD training of a *Perceptron*, on the data,  $\mathbf{x}^n$ , to learn the optimal weights,  $\mathbf{w}$ , as a Langevin process.<sup>20</sup> The power of the **StatMech** approach comes from the fact that the core

<sup>20</sup>Typically, we have no guarantee of the true equilibrium in a high-dim nonconvex landscape; so, when the *Thermodynamic limit* exists, the Langevin process converges or relaxes to the thermodynamic equilibrium.

Statistical Physics	Neural Network Learning
Gaussian field variables	Gaussian i.i.d data $\xi^N \in \mathcal{D}$
State Configuration	Trained / Learned weights $\mathbf{w}$
State Energy Difference	Training and Generalization Errors $\bar{\mathcal{E}}_{train}, \bar{\mathcal{E}}_{gen}$
Temperature	Amount of stochasticity present during training $T$
Annealed Approximation	Average over the data first
Thermal Average	Expectation w.r.t. the distribution of trained models
Free Energy	Generating function for the error(s) $F$

Table 2: Mapping between states and energies of a physical system and parameters of the learning process of a neural network.

SETOL Terminology	Explanation
Model Quality $\bar{\mathcal{Q}}$	Generalization accuracy, in the AA and at high-T
Layer Quality $\bar{\mathcal{Q}}$	Layer contribution to the accuracy, in the AA and at high-T
Layer Quality-Squared $\bar{\mathcal{Q}}^2$	Layer Quality squared, used for technical reasons
Quality Generating Function $\Gamma_{\bar{\mathcal{Q}}}, \Gamma_{\bar{\mathcal{Q}}^2}$	Generating function for quality
Annealed Hamiltonian $H^{an}$	Energy function, for errors or accuracies
Effective Hamiltonian $H^{eff}$	Exact energy function, but restricted to a low-rank subspace

Table 3: Additional terminology introduced for the SETOL. Notice that the Quality Generating Function  $\Gamma$  is simply one minus the Free Energy,  $\Gamma := 1 - F$ , but it introduced because sign convention for the Free Energy is always decreasing with the error. In contrast, we define the Hamiltonian in terms of the model error or accuracy, depending on the context.

concept of Thermal average corresponds to taking the expectation of a given quantity only *over the set of trained models*, as opposed to uniformly over all possible models (or, in a worst-case sense, over all possible models in a model class). This capability is particularly compelling in light of the StatMech capacity to characterize disordered systems with complex non-convex energy landscape (which can even be *glassy*, characterized by a highly non-convex landscape [?, ?, ?], and recognized classically by a slowing down of the training dynamics [?]). Thus, concepts such as training and Generalization Error arise naturally from integrals that are familiar to StatMech; and theoretical quantities such as Load, Temperature, and Free Energy also map onto useful and relevant concepts [?]. The Free Energy is of particular interest because it can be used as a generating function to obtain the desired Generalization Error and/or accuracy. We wish to understand how to compute the associated thermodynamic quantities such as the expected value of the various forms of the Average Generalization Error ( $\bar{\mathcal{E}}_{gen}$ ), Partition Function ( $Z$ ), and the Free Energy ( $F$ ) and other Generating Functions ( $\Gamma$ ).

**The Student-Teacher model.** We seek to compute and/or estimate the Average Generalization Accuracy for a *fixed* Teacher Perceptron  $T$  by averaging over an ensemble of Student  $S$  Perceptrons, in the Annealed Approximation (AA), and at High-Temperature (high-T); we call this ST Model Quality, and denote it  $\bar{\mathcal{Q}}^{ST}$ . We will then, later, in Section 5, we generalize  $\bar{\mathcal{Q}}^{ST}$  to an arbitrary

layer in Multi-Layer Perceptron, giving a Layer Quality, i.e.,  $\bar{Q}^{ST} \rightarrow \bar{Q}$ . This formulation of the ST problem is different than the classic approach because one usually does not fix the Teacher but, instead, averages over all Teacher vectors  $\mathbf{t}$  [?, ?]. This is one of many ways that distinguishes the current approach from previous work. Because of this, we present both a pedagogic derivation of  $\bar{Q}^{ST}$  (for a general NN in Section 4.2, and for the ST model specifically in the Appendix, Section A.2), and we provide a more heuristic approach in Section 4.3.2, assuming the AA and high-T at all times.

**Towards a Semi-Empirical Theory.** In the SETOL approach to StatMech, we want a matrix generalization of the ST Model Quality,  $\bar{Q}^{ST}$ , for a single Layer Quality  $\bar{Q} \sim \bar{Q}_L^{NN}$  in an arbitrary Multi-Layer Perceptron (MLP). This matrix generalization is a relatively straightforward extension of the classical (i.e., for a vector Teacher) SMOG ST Model Quality (but our SETOL approach will use it in a conceptually new way). In our matrix generalization, the Teacher vector  $\mathbf{t}$  is replaced by a Teacher matrix  $\mathbf{T}$  (i.e.,  $\mathbf{t} \rightarrow \mathbf{T}$ ); and, in our SETOL approach,  $\mathbf{T}$  will be an actual (pre-)trained NN weight matrix (i.e.,  $\mathbf{T} = \mathbf{W}$ ). Importantly, this matrix  $\mathbf{W}$  is neither a Gaussian Random Matrix, nor is it obtained from Gaussian i.i.d training data. As such, for our SETOL theory, we seek an expression that can be parameterized by the Teacher, and in particular by the ESD of the Teacher. This is what makes the basic method Semi-Empirical: even though we do not know the form of the Teacher, we make a modeling assumption that the Teacher has the same limiting spectral distribution as the Student, and hence the same PL fit parameter  $\alpha$ . This is all done with the understanding that later we will augment (and hopefully “correct”) our mathematical formulations with phenomenological parameters fit from experimental data. To make the Semi-Empirical method a Semi-Empirical *Theory*, we not only seek to parameterize our model; but we also try to understand how to derive HTSR empirical metrics, such as Alpha and AlphaHat, how they arise from this formalism, how they are related to the correlations in the data, and why they are transferable and exhibit Universality. This gives what we call a Semi-Empirical *Theory*.

## 4.2 Mathematical Preliminaries of Statistical Mechanics

**SubSection Roadmap** Briefly, in the following subsection, we start by defining an arbitrary NN model, with weights ( $\mathbf{w}$ ) Then, we explain the difference between using real-world ( $\mathbf{x}$ ) and random ( $\xi$ ) data This lets us define an energy error function,  $\Delta E_{\mathcal{L}}(\mathbf{w})$ , the error the NN makes on the data. We then explain how to take different kinds of *Thermodynamic* averages of the data, including *Sample* and *Thermal average* and the implications, and the difference between computing errors and accuracies. Next, we define the *FreeEnergy* ( $F$ ) for the error(s), and the *GeneratingFunction* ( $\Gamma$ ) for the accuracy and/or quality. From here, we explain the *AnnealedApproximation* (AA) and how to define the *Annealed Hamiltonian*,  $H^{an}(\mathbf{w})$ , a crucial expression that will be the starting point later for our matrix model. In the AA,  $H^{an}(\mathbf{w})$  simplifies to  $H_{hT}^{an} = \epsilon(\mathbf{w})$ , where  $\epsilon(\mathbf{w})$  is an Annealed Error Potential that depends only on the weights  $\mathbf{w}$ . Likewise, we can define the Self-Overlap,  $\eta(\mathbf{w}) := 1 - \epsilon(\mathbf{w})$ , which is useful for obtaining the Quality. We show how to obtain the *Average Training and Generalization Errors*  $\bar{\mathcal{E}}_{train}, \bar{\mathcal{E}}_{gen}$  using the StatMech formalism, which defines them in terms of partial derivatives of the Free Energy ( $F$ ). Doing this, we show that in the AA and at high-T they are equivalent,  $[\bar{\mathcal{E}}_{train}^{ST}]^{an,hT} = [\bar{\mathcal{E}}_{gen}^{ST}]^{an,hT}$ , and can both be expressed as a Thermal Average over all Students, as a function of the Teacher, as  $[\bar{\mathcal{E}}_{gen}^{ST}]^{an,hT} = \langle H_{hT}^{an}(R) \rangle_{\mathbf{s}}^{\beta} = \langle \epsilon(R) \rangle_{\mathbf{s}}^{\beta}$ . Note that these averages are obtained by using the Free Energy as a Generating Function. We then explain how to obtain the Model Quality as partial derivatives of a *Generating Function* ( $\Gamma_{\bar{Q}}$ ). We then discuss the more advanced techniques, the *Large-N Approximation* and the *SaddlePointApproximation* (SPA), which will be used extensively later.



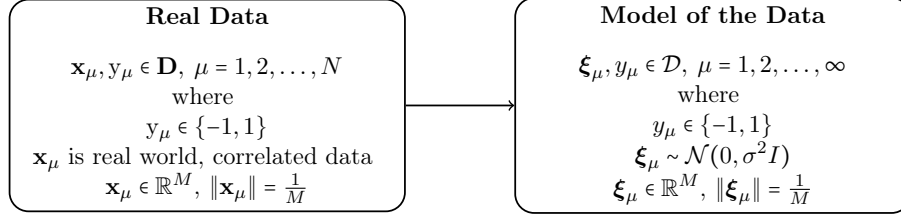


Figure 6: Mapping from a fixed set of  $n = N$  real-world, correlated data instances  $[\mathbf{x}, y] \in \mathbf{D}$  to an uncorrelated, random model of the data  $[\boldsymbol{\xi}, y] \in \mathcal{D}$ , drawn from a Gaussian i.i.d. distribution.

Finally, we introduce HCIZ integrals, which will be necessary to evaluate the matrix-generalized form of  $\Gamma_{\bar{\mathcal{Q}}}$  to obtain the final result.

In this subsection, we will compare and contrast several types of averages and energies we will encounter.

**4.2.1 Setup** In Section 4.2.1, we will start by describing the basic setup of the problem, including the distinction between the actual training process and how we model the training process.

**4.2.2 Sample Averages, Expected Values, and Thermal Averages** In Section 4.2.2, we will describe Thermal Averages (over the weights  $\mathbf{w}$ ) and Sample Averages (over the data  $\mathbf{x}$ )—in particular, under the Annealed Approximation (AA) and in the High-Temperature (High-T) limit—showing how they relate to each other and to the notion of Replica Averages.

**4.2.3 Free Energies and Generating Functions** In Section 4.2.3, we will make a connection between these different averaging notions and Free Energies and Generating Functions, showing how they relate to each other.

**4.2.4 The Annealed Approximation (AA) and the High-Temperature Approximation (high-T)** In Section 4.2.4, we explain the Annealed Approximation, the High-Temperature approximation, and the Thermodynamic Large-N limit and the Saddle Point Approximation (SPA). We also introduce the Quality Generating Function  $\Gamma_{\bar{\mathcal{Q}}}$

**4.2.5 Average Training and Generalization Errors and their Generating Functions** In Section 4.2.5, we will show how to compute the Average Training and Test/Generalization Errors  $\bar{\mathcal{E}}_{train}, \bar{\mathcal{E}}_{gen}$  using the Free Energy as a Generating Function, and how these errors are related to each other in the AA and High-T limit. We will also make connections with the Saddle Point Approximation (SPA).

**4.2.6 The Thermodynamic Large-N limit and the Saddle Point Approximation (SPA)** In Section 4.2.7, we discuss the large- $N$  Thermodynamic limit and the Saddle Point Approximation (SPA). We also introduce concept of Self-Averaging.

**4.2.7 HCIZ Integrals** Finally, in Section 4.2.8, we will describe how to express the Free Energy as a matrix-generalized Thermal Average over random matrices, called an HCIZ integral.

The various symbols and other important results are summarized in the Appendix A.1

#### 4.2.1 Setup

A basic issue in formulating SETOL is that one typically trains one large (expensive) NN, i.e., one does not split the data into training and testing sets. Thus, we want a methodology to

approximate quantities such as the generalization error or generalization accuracy that does not rely on traditional train-test splitting methods. To accomplish this, we will “model” the data, and we will use **StatMech** to construct quantities (basically, free energies or generating functions) so that we can compute training/testing errors by taking appropriate derivatives of these quantities.

In more detail, we imagine training a NN on  $n$  training data instances,  $\mathbf{x}_\mu$ , which are  $m$ -dimensional vectors, with labels  $y_\mu$ , chosen from a large but finite-size training data set  $\mathbf{D}$ . The goal of training the NN is to learn the  $m$  weights of the vector  $\mathbf{w}$  (or, later, a weight matrix  $\mathbf{W}$ ) by running gradient descent to minimize a loss function  $\mathcal{L}$  ( $\ell_2$ , cross-entropy, etc.). We want to model the real-world NN using a simple model from which we can obtain analytic expressions for the *Free Energy* and *Generating Function* we then use to compute Thermodynamic averages such as the *Average Generalization Error* ( $\bar{\mathcal{E}}_{gen}$ ) and *Model Quality* ( $\bar{Q}$ ) (which is our approximation to the *Average Generalization Accuracy*).

**Counting Samples and Features:  $n$ ,  $m$ ,  $N$ , and  $M$ .** We let the number of training samples be  $n$  and the dimension (i.e., number of features) for each sample be  $m$ . For simplicity, we also use  $N$  (instead of  $n$ ) and  $M$  (instead of  $m$ ), recognizing that *in later sections*  $N$  and  $M$  will refer to the dimensions of a layer’s weight matrix (i.e., an  $N \times M$  matrix). We stress that here, in this subsection,  $n = N$  and  $m = M$  only hold for our immediate analysis, to avoid extra notation. When we move to matrix-based analyses, we will revisit (and possibly distinguish)  $N$  (layer input dimension) and  $n$  (training-set size).

**Actual and Model Data and Energies.** Consider having a large set  $n = N$  of actual, real-world data,

$$\mathbf{x}_\mu, y_\mu \in \mathbf{D}, \mu = 1, \dots, n, \quad (17)$$

where  $\mathbf{x}_\mu$  is an  $m$ -dimensional real vector,  $\mathbf{x}_\mu \in \mathbb{R}^m$ ,  $y_\mu$  is a binary label taking values 1 or  $-1$ , denoted  $\{-1, 1\}$ , and  $\mathbf{D}$  denotes the finite-size dataset.

Notice that  $\mathbf{x}_\mu \in \mathbb{R}^m$  is normalized such that the Frobenius norm squared is  $\frac{1}{m}$ :

$$\|\mathbf{x}_\mu\|_F^2 := \sum_{i=1}^m \mathbf{x}_{\mu,i}^2 = \frac{1}{m} \quad (18)$$

We call  $\mathbf{x}^n$ , an  $n$ -dimensional sample (of the training data instances  $\mathbf{x}$ ) from  $\mathbf{D}$ . We may or may not specify the labels for this sample, depending on the context

We associate model errors with a (change in) energy  $\Delta E_{\mathcal{L}}$ . Smaller energies correspond to smaller errors and therefore better models. For example, for the mean-squared-error (MSE) loss, one has

$$\Delta E_{\mathcal{L}}(\mathbf{w}, \mathbf{x}_\mu, y_\mu) := (y_\mu - E_{NN}(\mathbf{w}, \mathbf{x}_\mu))^2, \quad (19)$$

where  $E_{NN}(\mathbf{w}, \mathbf{x}_\mu)$  is output prediction of the NN, as in Eqn. 1.

To estimate quantities such as the generalization error or generalization accuracy, we will adopt an approach that involves replacing the real data with Gaussian data and the NN with a parametric model that we will fit with a Semi-Empirical procedure (described later). To model the data, we specify the following mapping:

$$\mathbf{D} \rightarrow \mathcal{D}, \quad \mathbf{x}_\mu \rightarrow \boldsymbol{\xi}_\mu, \quad y_\mu \rightarrow y_\mu, \quad (20)$$

where we denote the model training and/or test data instances as  $(\boldsymbol{\xi}, y)$  such that

$$\boldsymbol{\xi}_\mu, y_\mu \in \mathcal{D}, \quad \mu = 1, \dots, \infty. \quad (21)$$



Here,  $\xi_\mu$  is a random vector (i.e., an  $m$ -dimensional random variable,  $\xi_\mu \in \mathbb{R}^m$ ), sampled from an i.i.d Gaussian distribution  $\mathcal{D}$ , and  $y_\mu$  denotes the (binary) label and/or NN output. We call  $\xi^n$  an  $n$ -dimensional *Model Sample* from  $\mathcal{D}$ .

#### 4.2.2 BraKets, Expected Values, and Thermal Averages

Given the setup from Section 4.2.1, we will want to model the average (change in) energy,  $\Delta E_{\mathcal{L}}(\mathbf{w}, \mathbf{x}^n)$ , averaged over some  $n$ -size data set  $\mathbf{x}^n$ . We can write the Total Data Sample Error as using an overloaded, operator notation

$$\Delta E_{\mathcal{L}}(\mathbf{w}, \mathbf{x}^n, \mathbf{y}^n) := \sum_{\mu=1}^n \Delta E_{\mathcal{L}}(\mathbf{w}, \mathbf{x}_\mu, y_\mu), \quad (22)$$

where the boldface  $\Delta E_{\mathcal{L}}(\mathbf{w})$  indicates this a sum over the entire set of  $n$  pairs  $[\mathbf{x}^n, \mathbf{y}^n]$ . We should keep in mind that this depends on the specific set of  $n$  data pairs  $[(\mathbf{x}_\mu, y_\mu) \in \mathbf{D} | \mu = 1, \dots, n]$ , although later we will model the labels  $y_\mu$  as the output of another NN when describing the Student-Teacher model. For that reason, for now, we will assume that the  $y_\mu$  is *implicit* in  $\Delta E_{\mathcal{L}}$ , will drop the  $y_\mu$  and  $\mathbf{y}^n$  symbols, and just write this total error / energy difference as

$$\Delta E_{\mathcal{L}}(\mathbf{w}, \mathbf{x}^n) := \sum_{\mu=1}^n \Delta E_{\mathcal{L}}(\mathbf{w}, \mathbf{x}_\mu), \quad (23)$$

which is now a function of the entire set of  $n$  vectors  $[\mathbf{x}^n]$  (where the labels  $\mathbf{y}$  have been set implicitly).<sup>21</sup> This operator notation will provide useful later in Section 4.3.2 (see Eqn. 85) and in Appendix A.2.

We will not, however, work directly with Samples and Sample Averages. Instead, we will model them. To that end, we need to estimate them with a theoretical approach. For example, we can write the Total Data Sample Error in terms of our random data variables  $\xi$  formally as

$$\Delta E_{\mathcal{L}}(\mathbf{w}, \xi^n) := \sum_{\mu=1}^n \Delta E_{\mathcal{L}}(\mathbf{w}, \xi_\mu), \quad (24)$$

but to evaluate this we need to take an integral and/or Expected Value over the data sample  $\xi^n$ .

**Expected Values.** We need to compute various sums and integrals, sampling from a model  $\mathcal{D}$  for the actual data distribution  $\mathbf{D}$ , which will frequently (but not always) be defined as more familiar Expected Values. We will denote Expected Values using physics Bra-Ket notion. Importantly, we use the term Expected Value in the physics sense, and BraKets will denote an un-normalized sum or integral; when the quantity is to be normalized, we will denote the normalization explicitly. For example, given a function  $f(\xi)$ , we write the BraKet integral as:

$$\langle f(\xi) \rangle_\xi := \int d\mu(\xi) f(\xi), \quad (25)$$

and if we want to express an  $n$ -dimensional average over  $f()$  then we would express this as

$$\langle f(\xi^n) \rangle_{\xi^n} := \frac{1}{n} \int d\mu(\xi^n) f(\xi^n) \quad (26)$$

<sup>21</sup>In the classic Student Teacher model, the labels  $\mathbf{y}^n$  represent the Teacher outputs and are effectively treated as uniform random variables to be averaged over later. In this work, the Teacher is fixed so we can drop the labels.

where the BraKet  $\langle \dots \rangle_{\xi^n}$  deotes the integral over the Model Data  $\xi^n$ , but with the convention that the normalization  $\frac{1}{n}$  appears inside the Bra-Ket implicitly. If this integral is normalized properly, then this denotes the familiar Expected Value  $\mathbb{E}_{\xi}[f(\xi)]$ . For a more complicated example, consider how to compute Expected Value of the Data Sample Error. That is, we want to model the average Data Sample Error using:

$$\frac{1}{n} \Delta \mathbf{E}_{\mathcal{L}}(\mathbf{w}, \mathbf{x}^n) \xrightarrow{\text{Expected Value}} \langle \Delta \mathbf{E}_{\mathcal{L}}(\mathbf{w}, \xi^n) \rangle_{\xi^n}, \quad (27)$$

In this case, we obtain:

$$\begin{aligned} \langle \Delta \mathbf{E}_{\mathcal{L}}(\mathbf{w}, \xi^n) \rangle_{\xi^n} &:= \frac{1}{n} \int d\mu(\xi^n) \Delta \mathbf{E}_{\mathcal{L}}(\mathbf{w}, \xi^n) \\ &= \frac{1}{n} \int \prod_{\mu=1}^n d\xi_{\mu} P(\xi^n) \sum_{\mu=1}^n \Delta E_{\mathcal{L}}(\mathbf{w}, \xi_{\mu}), \end{aligned} \quad (28)$$

where  $P(\xi^n)$  is an  $n$ -dimensional probability distribution (i.e., an  $n$ -dimensional Gaussian distribution), normalized to 1, and where the subscript  $\xi^n$  on the Ket reminds us this is an average or Expected Value of a finite,  $n$ -size Model Sample. (This is used in both Sections 4 and 5.) The normalization  $\frac{1}{n}$  is included in the defintion to ensure the Bra-Ket is a proper Expected Value. The measure  $d\mu(\xi^n)$  signifies a  $n$  i.i.d Gaussian vector  $\xi$ , drawn from an  $m$ -dimensional data model vectors  $\xi$ . Also, in some cases, we make use the subscript  $\xi^n$  on the Ket as  $\langle \dots \rangle_{\xi^n}$ ; this represents an integral over the data, but not an average or Expected Value.

**Size-Extensivity, Size-Intensivity, and Size-Consistency** A key requirement for the Thermodynamic limit in **StatMech** is *Size-Extensivity*: that physically meaningful quantities (i.e, total energies and free energies) scale linearly with the system size,  $n = N$ . Along with this, Thermodynamic average quantities should be *Size-Intensive*, meaning that they remain independent of  $n = N$  as the system size increases. In our setting, Size-Extensivity and Size-Intensivity underpin the large- $N$  limit, ensuring that macroscopic observables become independent of microscopic fluctuations.

As an example of Size-Extensivity and Size-Intensivity, we write the Expected Value (i.e., the data-average) of Data Sample Error  $\Delta \mathbf{E}_{\mathcal{L}}(\mathbf{w}, \xi^n)$  (Eqn. 28) in the large- $N$  limit as

$$\lim_{n \gg 1} \langle \Delta \mathbf{E}_{\mathcal{L}}(\mathbf{w}, \xi^n) \rangle_{\xi^n} = \lim_{n \gg 1} \frac{1}{n} \int \prod_{\mu=1}^n d\xi_{\mu} P(\xi^n) \sum_{\mu=1}^n \Delta E_{\mathcal{L}}(\mathbf{w}, \xi_{\mu}). \quad (29)$$

Here, the notation ( $n \gg 1$ ) means  $n$  grows arbitrarily large, but is not necessarily at the limit point ( $n = \infty$ ). The Total Data Sample Error  $\Delta \mathbf{E}_{\mathcal{L}}(\mathbf{w}, \xi^n)$  is Size-Extensive, whereas the average  $\langle \Delta \mathbf{E}_{\mathcal{L}}(\mathbf{w}, \xi^n) \rangle_{\xi^n}$  is Size-Intensive. This limit will be implicit later when taking a Saddle Point Approximation (see below).<sup>22</sup>

The data-averaged error  $\langle \Delta \mathbf{E}_{\mathcal{L}}(\mathbf{w}, \xi^n) \rangle_{\xi^n}$  will appear frequently below. For compatibility with [?], we denote it using the symbol  $\epsilon(\mathbf{w})$ :

$$\epsilon(\mathbf{w}) := \lim_{n \gg 1} \langle \Delta \mathbf{E}_{\mathcal{L}}(\mathbf{w}, \xi^n) \rangle_{\xi^n} \quad (\text{Size-Intensive}). \quad (30)$$

where, by our normalization here,  $\epsilon(\mathbf{w}) \in [0, 1]$ . The symbol  $\epsilon(\mathbf{w})$  is our theoretical estimate of the sample average  $\Delta \mathbf{E}_{\mathcal{L}}(\mathbf{w}, \xi^n)$  (Eqn. 31), well-defined for any  $n$ . Q: Is that eqn a typo, should be

<sup>22</sup>As we are working within a “physics-level of rigor, we take some liberties in evaluating these large- $N$  limits; and we leave the formal proofs for future work.

Eqn. 28? If yes, this sentence should be up there. We also call  $\epsilon(\mathbf{w})$  the *Annealed Error Potential*, which will be made clear below.

It is also convenient to write *Total Annealed Error Potential* as an Energy,

$$\mathcal{E}(\mathbf{w}) := n\epsilon(\mathbf{w}) \quad (\text{Size-Extensive}). \quad (31)$$

This will only be useful when the Thermodynamic limit exists, and this can be reasonably expected for the Annealed Approximation (AA), which is the regime in which SETOL will be developed.<sup>23</sup>

There is a related notion of *Size-Consistency*, often introduced in Quantum Chemistry through the Linked Cluster Theorem [?, ?], which states that the *average* energies and/or free energies ( $\bar{F}$ ) of in particular a correlated system scale with  $\mathcal{M}$ , the number of “independent interacting components”:

$$\begin{aligned} -\beta\bar{F} = \langle \ln Z \rangle_{\xi^n} &= \sum_{\mu=1}^{\mathcal{M}} (\text{Connected Components}) \\ &= \sum_{\mu=1}^{\mathcal{M}} \text{Cumulants}(\mu) + \text{higher order terms} \end{aligned} \quad (32)$$

For SETOL, below, these connected components will be matrix-generalized cumulants from RMT. For NNs, Size-Consistency appears when scaling the number of features in our matrix model, and it ensure that our layer and model Qualities remain well-behaved as we increase  $M$ , just as we do for  $N$ . For a simple example, see Appendix A.2 where we derive the expression for the matrix-generalized Annealed Hamiltonian  $H^{an}$ . Both Size-Extensivity (in  $N$ ) and Size-Consistency (in  $M$ ) are crucial in our SETOL analysis: they justify taking the large- $N$  approximation for matrix integrals, and they ensure our resulting the HCIZ integral—a sum of integrated RMT cumulants (below)—scales with the dimension of the Effective Correlation Space (ECS),  $\mathcal{M} = \tilde{M}$ .

**From Errors to Accuracies: The Average Generalization Accuracy, the Quality, and the Self-Overlap.** We have been discussing various forms of errors. In SETOL, we will, however, primarily be concerned with approximating the *Average Generalization Accuracy*, or, more generally, the Quality of a NN model and/or its layers.<sup>24</sup> The average accuracy is simply one minus the error. To represent this, we introduce the *Self-Overlap*  $\eta$ , which is defined generally as

$$\eta(\mathbf{w}) := 1 - \epsilon(\mathbf{w}) \in [0, 1], \quad (33)$$

and which describes the “overlap” between the true and the predicted labels. Unlike here, however, in later sections (4.3.2, 5.1, and Appendix A.3) we will first define a data-dependent Self-Overlap, so that we may obtain  $\eta(\mathbf{w}) := \langle \eta(\mathbf{w}, \xi) \rangle_{\xi^n}$  directly.

**Bracket Notation.** We will use physics Bra-Ket notation,  $\langle \dots \rangle$ , to denote different kinds of sums and integrals, with superscripts and subscripts, and for Expected Values (estimated theoretical averages). We use superscripts to denote the kind of integral or average:

Thermal  $\langle \dots \rangle^\beta$ , Annealed  $\langle \dots \rangle^{an}$ , high-T  $\langle \dots \rangle^{hT}$ , HCIZ  $\langle \dots \rangle^{IZ}$ , etc.

We use subscripts to emphasize the dependent variables:

<sup>23</sup>We should note that, while our model training and generalization errors are always expressed energies, an energy is not necessarily a model error.

<sup>24</sup>Technically, the Quality will estimate the average *Precision* rather than the Accuracy. This will distinction will be clarified in the Section 4.3.

1129

weights  $\langle \cdots \rangle_{\mathbf{w}}, \langle \cdots \rangle_{\mathbf{s}}, \langle \cdots \rangle_{\mathbf{s}}$ 

1130

data  $\langle \cdots \rangle_{\xi}, \langle \cdots \rangle_{\xi^n}, \langle \cdots \rangle_{\xi^N}, \langle \cdots \rangle_{\bar{\xi}^n}, \langle \cdots \rangle_{\bar{\xi}^N}$ 

1131

When averaging over the data, the subscript will appear with a bar (i.e.  $\bar{\xi}^n$ ), but when just integrating over the data, no bar will appear (i.e.,  $\xi^n$ ). We also reuse these symbols for other quantities, such as the  $Z_n^{an,hT}$ ,  $\bar{\mathcal{E}}_{gen}^{an,hT}$ ,  $H^{an}(\mathbf{w})$ , etc, but may mix-and-match subscripts and superscripts for visual clarity.

1135

**Sign Conventions.** Finally, we discuss the sign conventions used. Since errors decrease with better models, Energies ( $\Delta E_{\mathcal{L}}(\mathbf{w}, \xi), \mathcal{E}(\mathbf{w}), \epsilon(\mathbf{w}), \cdots$ ) and Free Energies ( $F$ ) are minimized to obtain better models. Likewise, since accuracies increase with better models, Qualities ( $\bar{Q}, \bar{Q}^2, \cdots$ ), Self-Overlap ( $\eta$ ), and Quality Generating Function ( $\Gamma$ ) would be maximized to obtain better models. An exception will be Hamiltonians ( $H, \mathbf{H}$ ), which will depend on context.

1140

**Thermal Averages (over weights).** To evaluate the expectation value of some equilibrium quantity that depends on the weights  $\mathbf{w}$  (say  $\mathbb{E}_{\mathbf{w}}^{\beta}[f(\mathbf{w})]$ ), one uses a Thermal Average. By this, we mean a *Boltzmann-weighted average*: given a function  $f(\mathbf{w})$ , we define the Thermal Average over  $\mathbf{w}$  as

1143

$$\langle f(\mathbf{w}) \rangle_{\mathbf{w}}^{\beta} := \frac{1}{Z_n} \int d\mu(\mathbf{w}) f(\mathbf{w}) e^{-\beta \mathcal{E}(\mathbf{w})}, \quad (34)$$

1144

where the superscript  $\beta$  denotes Thermal Average,  $\beta = \frac{1}{T}$  is an inverse temperature, and  $Z_n$  is the normalization term (or Partition function), defined as

1145

$$Z_n := \int d\mu(\mathbf{w}) e^{-\beta \mathcal{E}(\mathbf{w})}, \quad (35)$$

1146

defined for the  $n$ -size *Data Sample*  $[\xi^n]$ . In particular, when we want to compute the Thermal Average of the *Total Energy* difference or Error  $\mathcal{E}(\mathbf{w})$  over  $\mathbf{w}$ , we could write

1147

$$\langle \mathcal{E}(\mathbf{w}) \rangle_{\mathbf{w}}^{\beta} := \frac{1}{Z_n} \int d\mu(\mathbf{w}) \mathcal{E}(\mathbf{w}) e^{-\beta \mathcal{E}(\mathbf{w})}. \quad (36)$$

1148

Importantly, we will never calculate the average errors directly like this. Instead, we will calculate them from partial derivatives of the Free Energy  $F$  (as shown below). Also, in some cases, we may use  $\langle \cdots \rangle_{\xi^n}^{\beta}$  to denote what looks like a Thermal Average over the data; this will be explained below when it occurs.

1151

1152

**Other Notation: Overbars, Superscripts and Subscripts.** As above, we may also occasionally denoted averages using the common notation for expected values,  $\mathbb{E}[\cdots]$ . See Table 7 and 8 in Appendix A for a list of these and other notational conventions and symbols we use.

1155

When discussing quantities such as the Free Energy ( $F$ ), training and test errors/energies ( $\mathcal{E}$ ), the Layer Quality ( $\bar{Q}$ ), etc., we will place a bar over the symbol (i.e.,  $\bar{F}$ ,  $\bar{\mathcal{E}}$ ,  $\bar{Q}$ , etc.) when referring to an average over  $n$  (or  $N$ , below). Otherwise, we will refer to these quantities as the total (averaged) energy, error, quality, etc.

1158

1159

Finally, when referring to the model (i.e., theoretical) training and generalization errors, we will use the superscript  $ST$  for the average Student-Teacher training and generalization errors,  $\bar{\mathcal{E}}_{train}^{ST}$  and  $\bar{\mathcal{E}}_{gen}^{ST}$ , respectively, and the superscript  $NN$  for the matrix-generalized NN layer average training and generalization errors,  $\bar{\mathcal{E}}_{train}^{NN}$  and  $\bar{\mathcal{E}}_{gen}^{NN}$ , respectively. When referring to empirical errors, we denoted these as  $\bar{\mathcal{E}}_{train}^{emp}$  and  $\bar{\mathcal{E}}_{gen}^{emp}$ , respectively.

1163

### 4.2.3 Free Energies and Generating Functions

If one needs an average energy (or error), it is often easier to calculate the associated Free Energy and take corresponding partial derivatives than it is to compute that quantity directly via an expected value or Thermal Average. Generally speaking, a Free Energy,  $F_n$ , is defined in terms of a partition function  $Z_n$  as

$$\beta F_n := -\ln Z_n. \quad (37)$$

Keep in mind that  $Z$  may actually be a function of the data  $\xi$  (or some other variables), i.e.,  $Z(\xi)$ , but we usually don't write this explicitly. Likewise, while both  $F_n$  and  $Z_n$  depend explicitly on the system size  $n$ , we will only include these subscripts when emphasizing this. Also,  $F$  has units of Energy or Temperature, so  $\beta F = -\ln Z$  is a kind-of unitless Free Energy. Each model (in single-layer models) and/or layer (in multi-layer models) will have its own Partition Function and associated Generating Functions. We call  $F$  and  $Z$  *Generating Functions* because they can be used to generate the model errors. That is, given an  $F$  and/or  $Z$ , we can “generate” the training and generalization errors with the appropriate partial derivatives [?, ?].

From this generating function perspective, i.e., when using a generating function to compute quantities of interest, we can work with other transformations of  $F$ . Most notably, we will consider

$$\Gamma = n - F. \quad (38)$$

where  $n$  is the number of free parameters ( $n = N$  for a vector model, but a matrix model has  $N \times M$  free parameters; see Appendix A.2.1). The quantity  $\Gamma$  decreases as the error increases (as opposed to  $F$ , which increases), i.e., it increases as the accuracy of quality of the model increases. Thus, we will use it as a generating function for the model quality/accuracy. For average Qualities, one has

$$\bar{\Gamma} := 1 - \bar{F} \quad (39)$$

for the model or layer under consideration (see below).

### 4.2.4 The Annealed Approximation (AA) and the High-Temperature Approximation (high-T)

In the traditional SMOG approach, one models the (typical) generalization behavior of a NN by defining and computing the Expected Value of the Free Energy of the model. The full expected value of the Free Energy,  $\beta F_n = -\ln Z_n$ , with respect to the (model) data  $\xi^n$ , is:

$$\mathbb{E}_{\xi^n}[\beta F_n] = \beta \bar{F} := -\langle \ln Z_n \rangle_{\xi^n}, \quad (40)$$

where  $\langle \ln Z_n \rangle_{\xi^n}$  means where we average over the  $n$  samples of the data ( $\xi^n$ , of size  $n$ ). This is also called the Quenched Free Energy. This is, however, frequently too difficult to analyze, and doing so typically requires a so-called Replica calculation.

The Annealed Approximation (AA) is a way of taking the data-average *first* and greatly simplifies the model under study and its analysis. The standard way to move forward is to follow the methods used in disordered systems theory [?, ?]. The mapping is:

Average over the Data	$\leftrightarrow$	Annealed Approximation	$\leftrightarrow$	Disorder Average
Learning the Weights	$\leftrightarrow$	NN Optimization Algorithm	$\leftrightarrow$	Thermal Average.

1196 **The Annealed Approximation (AA).** Formally, the AA makes the substitution

$$-\langle \ln Z_n \rangle_{\xi^n} \approx -\ln \langle Z_n \rangle_{\xi^n}. \quad (41)$$

1197 Here, we are *averaging over the disorder*. We may associate:

$$\begin{aligned} -\langle \ln Z_w(\xi^n) \rangle_{\xi^n} &: \text{the (unitless) Quenched Free Energy} \\ -\ln \langle Z_n(\xi^n) \rangle_{\xi^n} &: \text{the (unitless) Annealed Free Energy.} \end{aligned}$$

1198 Applying the AA amounts to applying Jensens inequality *as an equality*, and it allows will let us  
1199 interchange integrals and logarithms when computing the data average:

$$\ln \frac{1}{n} \int d\mu(\xi^n)(\dots) \xrightarrow{\text{AA}} \frac{1}{n} \int d\mu(\xi^n) \ln(\dots). \quad (42)$$

1200 This will allow us to switch the order of the data and the thermal averages, i.e.,

$$\langle \langle \dots \rangle_{\mathbf{w}} \rangle_{\xi^n}^{\beta} \xrightarrow{\text{AA}} \langle \langle \dots \rangle_{\xi^n}^{\beta} \rangle_{\mathbf{w}}, \quad (43)$$

1201 greatly simplifying the analysis.

1202 The use of the AA is common in **StatMech**, as it simplifies computations considerably; and it  
1203 is chosen when it holds exactly (if, say,  $x$  is a typical sample from  $\mathcal{D}$  and  $Z_w(\xi)$  has a well-defined  
1204 mean). In contrast, there are situations in **StatMech** when the average is atypical, and then it  
1205 one can get different results for the Quenched versus Annealed cases. In a practical sense, one  
1206 imagines this may occur when the data is very noisy and/or mislabeled, and this requires a special  
1207 treatment [?].

1208 **Annealed Hamiltonian  $H^{an}(\mathbf{w})$  and Annealed Partition Function  $Z^{an}$ .** When we apply  
1209 the AA (as in Eqn. 42), we average over the data  $\xi^n$  first. Doing this will allow us to develop a  
1210 theory in terms a (Temperature dependent) Annealed Error Potential.

1211 Following [?] (see their Eqn.(2.30)), we will call this average the *Annealed Hamiltonian*,  $H^{an}(\mathbf{w})$ ,  
1212 which we define as

$$\beta H^{an}(\mathbf{w}) := -\frac{1}{n} \ln \int d\mu(\xi^n) e^{-\beta \Delta \mathbf{E}_{\mathcal{L}}(\mathbf{w}, \xi^n)}. \quad (44)$$

1213 The Annealed Hamiltonian is a simple “mean-field-like” Hamiltonian for the problem. This can be  
1214 seen by noting that we can express Eqn. 44 as an integral over a single data example  $\xi$ :

$$\begin{aligned} \beta H^{an}(\mathbf{w}) &= -\frac{1}{n} \ln \int d\mu(\xi^n) e^{-\beta \sum_{\mu=1}^n \Delta E_{\mathcal{L}}(\mathbf{w}, \xi_{\mu})} \\ &= -\ln \left[ \int d\mu(\xi^n) e^{-\beta \sum_{\mu=1}^n \Delta E_{\mathcal{L}}(\mathbf{w}, \xi_{\mu})} \right]^{\frac{1}{n}} \\ &= -\ln \left[ \prod_{\mu=1}^n \int d\mu(\xi_{\mu}) e^{-\beta \Delta E_{\mathcal{L}}(\mathbf{w}, \xi_{\mu})} \right]^{\frac{1}{n}} \\ &= -\ln \int d\mu(\xi) e^{-\beta \Delta E(\mathbf{w}, \xi)} \end{aligned} \quad (45)$$

1215 This will be a critical piece needed to generalize the vector-based ST Perceptron model to the  
1216 matrix-generalized ST MLP model. In BraKet notation, Eqn. 45 can be expressed as

$$\beta H^{an}(\mathbf{w}) := -\frac{1}{n} \ln \langle e^{-\beta \Delta \mathbf{E}_{\mathcal{L}}(\mathbf{w}, \xi^n)} \rangle_{\xi^n} = -\ln \langle e^{-\beta \Delta E(\mathbf{w}, \xi)} \rangle_{\xi}.$$



1217 Using  $H^{an}(\mathbf{w})$ , we can define the *Annealed Partition Function*,  $Z_n^{an}$ , as

$$Z_n^{an} := \int d\mu(\mathbf{w}) e^{-n\beta H^{an}(\mathbf{w})} \quad (46)$$

$$\begin{aligned} &= \int d\mu(\mathbf{w}) \exp \left[ -n \left( -\frac{1}{n} \ln \int d\mu(\xi^n) e^{-\beta \Delta \mathbf{E}_{\mathcal{L}}(\mathbf{w}, \xi^n)} \right) \right] \\ &= \int d\mu(\mathbf{w}) \exp \left[ \ln \int d\mu(\xi^n) e^{-\beta \Delta \mathbf{E}_{\mathcal{L}}(\mathbf{w}, \xi^n)} \right] \\ &= \int d\mu(\mathbf{w}) \int d\mu(\xi^n) e^{-\beta \Delta \mathbf{E}_{\mathcal{L}}(\mathbf{w}, \xi^n)}. \end{aligned} \quad (47)$$

1218 where the lines after the first line follow by substituting Eqn. 44 into Eqn. 46. Note also that the  
 1219 order of the integrals in Eqn. 47 is exactly what we expect using the AA, as in Eqn. 43. Also,  
 1220 analogously to Eqn. 45, we can write  $Z_n^{an}$  as the product of the  $n = 1$  case,  $Z_n^{an} = [Z_1^{an}]^n$ . Finally,  
 1221 we will only need the high-T version,  $H_{hT}^{an}$ , of  $H^{an}(\mathbf{w})$ , and this will take a very simple form. <sup>25</sup>

1222 **The High-Temperature (High-T) Annealed Hamiltonian** ( $H_{hT}^{an}(\mathbf{w}) = \epsilon(\mathbf{w})$ ) **and Partition**  
 1223 **Function** ( $Z_n^{an, hT}$ ). In addition to the AA, we will be evaluating our models at at high-T.  
 1224 Notably, the Annealed Hamiltonian  $H^{an}(\mathbf{w})$  in Eqn. 44 is a non-linear function of  $\beta$ ; by taking  
 1225 the high-T approximation, we can remove this dependence and obtain the simpler expression that  
 1226  $H^{an}(\mathbf{w}) = \epsilon(\mathbf{w})$ . This greatly simplifies both the Partition Function, i.e.,  $Z_n^{an, hT}$ , and subsequent  
 1227 results (below).

1228 To obtain the high-T result, we can write the Taylor expansion for  $e^{\beta \Delta \mathbf{E}_{\mathcal{L}}(\mathbf{w}, \xi^n)}$  and keep the  
 1229 first two terms:

$$e^{-n\beta \Delta \mathbf{E}_{\mathcal{L}}(\mathbf{w}, \xi^n)} \simeq 1 - \underbrace{\beta \Delta \mathbf{E}_{\mathcal{L}}(\mathbf{w}, \xi^n)}_{high-T} + (\beta \Delta \mathbf{E}_{\mathcal{L}}(\mathbf{w}, \xi^n))^2 + \dots \quad (48)$$

1230 Let us now express  $H^{an}(\mathbf{w})$  directly in terms of  $\epsilon(\mathbf{w}) = \langle \Delta \mathbf{E}_{\mathcal{L}}(\mathbf{w}, \xi^n) \rangle_{\xi^n}$  (see Eqn. 30) as a  
 1231 Thermal Average at high-T. To do so, let's take a high-T expansion of Eqn. 44 by expanding the  
 1232 exponential to first order in  $\beta$ , to obtain

$$\begin{aligned} \beta H_{hT}^{an}(\mathbf{w}) &= -\frac{1}{n} \ln \int d\mu(\xi^n) [1 - \beta \Delta \mathbf{E}_{\mathcal{L}}(\mathbf{w}, \xi^n)] \\ &\approx -\frac{1}{n} \int d\mu(\xi^n) \beta \Delta \mathbf{E}_{\mathcal{L}}(\mathbf{w}, \xi^n) \\ &= \langle \beta \Delta \mathbf{E}_{\mathcal{L}}(\mathbf{w}, \xi^n) \rangle_{\xi^n} \\ &= \beta \langle \Delta \mathbf{E}_{\mathcal{L}}(\mathbf{w}, \xi^n) \rangle_{\xi^n} \\ &= \beta \epsilon(\mathbf{w}). \end{aligned} \quad (49)$$

1233 Here, we have used the AA, the property that  $\ln(1 + y) \approx y$ , for  $|y| \ll 1$ , and the fact that  $\epsilon(\mathbf{w})$   
 1234 takes the form given in Eqn. 30. This gives  $H_{hT}^{an}(\mathbf{w}) := \epsilon(\mathbf{w})$ , which is no longer a non-linear  
 1235 function of  $\beta$ . Moving forward, we will assume we are taking the high-T limit like this.

<sup>25</sup>We will derive expressions for  $\epsilon(R)$  and  $\bar{\mathcal{E}}_{gen}^{ST}(R)$  in Eqn. 96 and Eqn. 97, respectively, using relatively simple arguments. In Appendix A.2.1 we present a more detailed derivation of  $H^{an}(R)$  and  $H_{hT}^{an}(R) = \epsilon(R)$  in Appendix A.2.2 we show that this derivations generalizes to the matrix-generalized case,  $H^{an}(\mathbf{R})$ . This more detailed derivation is important for our SETOL setup because it lets us define the normalization necessary for the TRACE-LOG condition.

Given Eqn. 50, we can now express the Annealed Partition Function at high-T directly in terms of the Annealed (i.e., data-averaged) error  $\epsilon(\mathbf{w})$ :

$$\begin{aligned} Z_n^{an,hT} &:= \int d\mu(\mathbf{w}) e^{-n\beta H_{hT}^{an}(\mathbf{w})} \\ &= \int d\mu(\mathbf{w}) e^{-n\beta \epsilon(\mathbf{w})} \end{aligned} \quad (51)$$

If we assume that at high-T we can make this approximation, since we only care about the small  $\beta$  results. This will prove very useful later when working with HCIZ integrals.

#### 4.2.5 Average Training and Generalization Errors and their Generating Functions

Here, we show how to derive the Average Training Error  $\bar{\mathcal{E}}_{train}$  and the Average Generalization Error  $\bar{\mathcal{E}}_{gen}$ , in the Annealed Approximation (AA), and at high-Temperature (high-T), using the Free Energy  $F$  and/or the Partition Function  $Z$  as a generating function. In particular, we show that, in the AA and at high-T, these errors are (approximately) equal and equal to the Thermal Average of the Annealed Error Potential,  $\bar{\mathcal{E}}_{train}^{an,hT} \approx \bar{\mathcal{E}}_{gen}^{an,hT} \approx \langle \epsilon(\mathbf{w}) \rangle_{\mathbf{w}}^\beta$ .

**Generating Functions for the Errors: the StatMech way.** In our theory, after applying the AA, we obtain expressions where the random model data  $\xi^n$  has been integrated out. This leaves formal quantities that depend only on the weights  $\mathbf{w}$ , which are the variables being learned during training. Since we are left with a distribution over  $\mathbf{w}$ , we define the training error not explicitly as an average over the training data itself, but instead in terms of how the Free Energy,  $\beta F := -\ln Z_n$ , varies with  $\beta$ , i.e., the amount of stochasticity in the model weights.

Following [?, ?], we define the *Average Training Error*, in the AA, by differentiating  $\ln Z_n^{an}$  with respect to  $\beta$ :

$$\bar{\mathcal{E}}_{train}^{an} := -\frac{1}{n} \frac{\partial(\ln Z_n^{an})}{\partial \beta} = -\frac{1}{n} \frac{1}{Z_n} \frac{\partial Z_n^{an}}{\partial \beta}. \quad (52)$$

This error captures how the model predictions will vary with changes in the learned weights  $\mathbf{w}$ , which implicitly describes how the changes will vary with the training data  $\xi^n$ . Similarly, we define the *Average Generalization Error*, in the AA, by differentiating  $\ln Z_n^{an}$  with respect to  $n$ , the number of data points. Following

$$\bar{\mathcal{E}}_{gen}^{an} := -\frac{1}{\beta} \frac{\partial(\ln Z_n^{an})}{\partial n} + \frac{1}{\beta} \ln z(\beta) = -\frac{1}{\beta} \frac{1}{Z_n} \frac{\partial Z_n^{an}}{\partial n} + \frac{1}{\beta} \ln z(\beta), \quad (53)$$

where  $z(\beta)$  is a constant normalization term that depends only on  $\beta$  (which, moving forward, we ignore, as it only shifts the scale). This error captures how the model's predictions will change as more data points are introduced.

In the Thermodynamic limit ( $n \gg 1$ ), these two definitions of the error become equivalent at High-T, and they equal to the Thermal Average of the Annealed Error Potential:

$$\bar{\mathcal{E}}_{train}^{an,hT} = \bar{\mathcal{E}}_{gen}^{an,hT} = \langle \epsilon(\mathbf{w}) \rangle_{\mathbf{w}}^\beta, n \gg 1. \quad (54)$$



1264 To see this, substitute Eqn. 51 into Eqn. 52, and take the partial derivative w.r.t  $\beta$ , to obtain

$$\begin{aligned}
\bar{\mathcal{E}}_{gen}^{an,hT} &:= \frac{1}{n} \frac{\partial(-\ln Z_n^{an,hT})}{\partial\beta} \\
&= -\frac{1}{n} \frac{\partial}{\partial\beta} \ln \int d\mu(\mathbf{w}) e^{-\beta n \epsilon(\mathbf{w})} \\
&= \frac{\frac{1}{n} \int d\mu(\mathbf{w}) n \epsilon(\mathbf{w}) e^{-\beta n \epsilon(\mathbf{w})}}{\int d\mu(\mathbf{w}) e^{-\beta n \epsilon(\mathbf{w})}} \\
&= \langle \epsilon(\mathbf{w}) \rangle_{\mathbf{w}}^\beta.
\end{aligned} \tag{55}$$

1265 Likewise, if we substitute Eqn. 51 into Eqn. 53, and take the partial derivative w.r.t  $n$ , we obtain

$$\begin{aligned}
\bar{\mathcal{E}}_{gen}^{an,hT} &:= \frac{1}{\beta} \frac{\partial(-\ln Z_n^{an,hT})}{\partial n} \\
&= -\frac{1}{n} \frac{\partial}{\partial\beta} \ln \int d\mu(\mathbf{w}) e^{-\beta n \epsilon(\mathbf{w})} \\
&= \frac{\frac{1}{\beta} \int d\mu(\mathbf{w}) n \epsilon(\mathbf{w}) e^{-\beta n \epsilon(\mathbf{w})}}{\int d\mu(\mathbf{w}) e^{-\beta n \epsilon(\mathbf{w})}} \\
&= \langle \epsilon(\mathbf{w}) \rangle_{\mathbf{w}}^\beta.
\end{aligned} \tag{56}$$

1266 Notice that both of these results arise because of the simple expression that appears in the exponent  
1267 of  $Z_n^{an,hT}$ , namely because  $-n\beta H_{hT}^{an}(\mathbf{w}) = n\beta \epsilon(\mathbf{w})$ .

1268 This equivalence reflects the fact that when the system is large enough, adding a new data  
1269 example to the training distribution is formally equivalent to adding noise, making the two errors  
1270 indistinguishable. This approach allows us to define both training and generalization errors in  
1271 terms of fundamental thermodynamic quantities, providing a simplified formal framework suitable  
1272 for empirical adjustment later.

1273 Also, note that the model data variables  $\xi$  do not enter the calculation because we integrated  
1274 them out before the calculation of Thermal Average. (This illustrates the difference between taking  
1275 an annealed versus a quenched average.) Also, our sign convention is consistent with a model of  
1276 NN training that *minimizes* the total loss ( $\mathcal{L}$ ) and/or ST error, and, therefore minimizes Free  
1277 Energies as well.

1278 More generally, we see that we can use the Partition Function,  $Z$ , and/or the Free Energy,  
1279  $\beta F := -\ln Z$ , as a Generating Function to obtain any desired Thermodynamic average by taking  
1280 the appropriate partial derivative of the corresponding form of  $\ln Z$ . In Appendix ?? we show how  
1281 to obtain  $\bar{\mathcal{E}}_{train}^{an}$  and  $\bar{\mathcal{E}}_{gen}^{an}$  obtain explicitly in this way using  $Z^{an}$ .

## 1282 4.2.6 The Quality ( $\bar{\mathcal{Q}}$ ) and its Generating Function ( $\Gamma_{\bar{\mathcal{Q}}}$ )

1283 Here, we explain how to define what we call the Quality  $\bar{\mathcal{Q}}$ , which is defined as the Thermal  
1284 Average of the Self-Overlap,  $\bar{\mathcal{Q}} := \langle \eta(\mathbf{w}) \rangle_{\mathbf{w}}^\beta$ , and which can be obtained from the associated  
1285 *Quality Generating Function*  $\Gamma_{\bar{\mathcal{Q}}}$ .

1286 For our purposes below, we define the Model Quality (as in Eqn. 7) as our approximation to  
1287 the Average Generalization Accuracy for our model. We denote the Model Quality for the ST  
1288 Perceptron model,  $\bar{\mathcal{Q}}^{ST}$ , and for a general NN,  $\bar{\mathcal{Q}}^{NN}$ , such that

$$\bar{\mathcal{Q}}^{ST} := 1 - \bar{\mathcal{E}}_{gen}^{ST} \tag{57}$$

$$\bar{\mathcal{Q}}^{NN} := 1 - \bar{\mathcal{E}}_{gen}^{NN} \tag{58}$$

1289 In this work, however, the Quality will always be defined at high-T, and so we may write

$$\bar{Q}^{ST} = 1 - [\bar{\mathcal{E}}_{train}^{ST}]^{an,hT} = 1 - [\bar{\mathcal{E}}_{gen}^{ST}]^{an,hT} \quad (59)$$

$$\bar{Q}^{NN} = 1 - [\bar{\mathcal{E}}_{train}^{NN}]^{an,hT} = 1 - [\bar{\mathcal{E}}_{gen}^{NN}]^{an,hT} \quad (60)$$

1290 We also define a Layer Quality, simply denoted  $\bar{Q}$ , which will describe the contributions an  
1291 individual layer makes to the overall Model Quality  $\bar{Q}^{NN}$ . To obtain the Layer Quality, we define  
1292 an accuracy-or Quality-Generating Function, denoted  $\Gamma_{\bar{Q}}$ , which is analogous to a layer Free  
1293 Energy, but with the opposite sign convention.

1294 Generally speaking, the Quality Generating Function  $\Gamma_{\bar{Q}}$  is defined in the AA, and at high-T  
1295 and is given as

$$\beta\Gamma_{\bar{Q}} := \ln \int d\mu(\mathbf{w}) e^{n\beta\eta(\mathbf{w})} \quad (61)$$

1296 For example, for the single-layer ST Perceptron,  $\Gamma_{\bar{Q}^{ST}} := n - F^{ST}$  (where  $n$  here is also the  
1297 number of free parameters for this model, and is in units of energy or error). The term  $\Gamma_{\bar{Q}^{ST}}$  is  
1298 directly related to the *Total* Free Energy  $F^{ST}$ , which can be seen by substituting Eqn. 33 for  $\eta(\mathbf{w})$   
1299 in Eqn. 61:

$$\begin{aligned} \beta\Gamma_{\bar{Q}^{ST}} &= \ln \int d\mu(\mathbf{w}) e^{n\beta(1-\epsilon(\mathbf{w}))} \\ &= \ln \int d\mu(\mathbf{w}) e^{n\beta} e^{-n\beta\epsilon(\mathbf{w})} \\ &= \ln \left( \int d\mu(\mathbf{w}) e^{n\beta} \right) + \ln \left( \int d\mu(\mathbf{w}) e^{-n\beta\epsilon(\mathbf{w})} \right) \\ &= \ln e^{n\beta} + \ln \left( \int d\mu(\mathbf{w}) e^{-n\beta\epsilon(\mathbf{w})} \right) \\ &= \beta n - \beta F^{ST} \end{aligned} \quad (62)$$

1300 Dividing by  $n$ , we can also recover the more general relation for the *Average* Free Energy,  $\bar{\Gamma}_{\bar{Q}} = 1 - \bar{F}$ .  
1301 (Eqn. 39). For the matrix case we do something similar; see Appendix A.3.

1302 Likewise, one can show that the Quality  $\bar{Q}$  (again, always in the AA and at high-T) can be  
1303 identified as the Thermal Average of the (data-averaged or Annealed) Self-Overlap,

$$\bar{Q} = \left\langle \langle \eta(\mathbf{w}, \xi) \rangle_{\xi^n} \right\rangle_{\mathbf{w}}^\beta = \langle \eta(\mathbf{w}) \rangle_{\mathbf{w}}^\beta \quad (63)$$

1304 We can then obtain  $\bar{Q}$  by taking the appropriate partial derivative of its Generating Function,  $\Gamma_{\bar{Q}}$ .

1305 For technical reasons, however, we will actually define and use a Generating Function for the  
1306 Average Layer Quality-Squared  $\bar{Q}^2$ , denoted  $\beta\Gamma_{\bar{Q}^2}^{IZ}$ . In analogy with Eqns. 52 and 53, and at  
1307 high-T and large-N (explained below), we can obtain  $\bar{Q}^2$  (see Section 5 Eqn. 127) as

$$\bar{Q}^2 := \lim_{N \gg 1} \frac{1}{\beta} \frac{\partial}{\partial N} \beta\Gamma_{\bar{Q}^2}^{IZ} \approx_{\text{high-T}} \lim_{N \gg 1} \frac{1}{N} \frac{\partial}{\partial \beta} \beta\Gamma_{\bar{Q}^2}^{IZ} \quad (64)$$

1308 See Section 5 and Appendix A.3 for more details.

#### 1309 4.2.7 The Thermodynamic Large-N limit and the Saddle Point Approximation (SPA)

1310 To evaluate  $\bar{Q}^2$ , we take a Large-N limit, which is in the Thermodynamic limit, which we assume  
1311 going forward.

1312 We note that, technically, the Thermodynamic limit and the Large-N limit are different. In  
 1313 particular, the Thermodynamic limit typically refers to the case where all Thermodynamic averages  
 1314 remain Size-Intensive as the system size increases; and, while this does refer to  $n = N$  growing  
 1315 large, frequently for NNs there is an additional constraint that the ratio ( $m/n = M/N$ ), i.e., the  
 1316 load, remains constant [?, ?]. For the ST model of the Perceptron, however, the Thermodynamic  
 1317 averages we need to compute not depend on  $m$ , so this constraint is not necessary. Moreover, later,  
 1318 when we form our matrix generalization of the ST model, we will not enforce that  $N/M$  remain  
 1319 constant but, instead, the final result will be Size-Consistent.

1320 Here, the Thermodynamic limit is the large-N limit of the model, as  $n = N$  grows arbitrarily  
 1321 large, i.e.  $n \gg 1$ . To express the average Free Energy  $\bar{F}$  in the large-N limit, we can write

$$-\beta\bar{F} = \lim_{n \gg 1} \frac{1}{n} \ln \int d\mu(\mathbf{w}) e^{-n\beta\epsilon(\mathbf{w})}. \quad (65)$$

1322 When this large-N approximation is well behaved, then the total energy  $\mathcal{E}(\mathbf{w})$  is extensive, i.e.,  
 1323 when  $\mathcal{E}(\mathbf{w}) = n\epsilon(\mathbf{w})$ ; and, consequently, the total Free Energy is also extensive, i.e.,  $F = n\bar{F}$ . (This  
 1324 is a cornerstone of statistical mechanics, as it allows for meaningful macroscopic predictions from  
 1325 microscopic interactions.)

1326 **Self-Averaging.** The existence of the limit signifies that the system is *Self-Averaging*, meaning  
 1327 that the macroscopic properties are independent of fluctuations, etc. This also implies that the  
 1328 relevant averages (i.e., training and generalization errors) are the same for almost all samples, or  
 1329 *realizations of the disorder*. Additionally, the Annealed and the Quenched averages,  $\ln\langle Z_n \rangle_{\xi^n}$  and  
 1330  $\langle \ln Z_n \rangle_{\xi^n}$ , respectively, become sharply peaked, and

$$\langle \ln Z_n \rangle_{\xi^n} \approx \ln\langle Z_n \rangle_{\xi^n}, \quad \text{as } n \rightarrow \infty. \quad (66)$$

1331 For a NN, Self-Averaging implies that the weights  $\mathbf{w}$  are *typical* of the distribution, and therefore  
 1332 the NN can generalize to similar but unseen test examples.

1333 **Saddle Point Approximation (SPA).** When the Thermodynamic limit exists, we can ap-  
 1334 proximate the asymptotic behavior of integrals we will encounter (e.g., the Free Energy  $F$  and/or  
 1335 Partition Function  $Z$ ) in the large-N limit by using the Saddle Point Approximation (SPA). In the  
 1336 case of the  $F$  and/or  $Z$ , we assume that we can apply the SPA:

$$\int d\mu(\mathbf{w}) e^{-n\beta\epsilon(\mathbf{w})} \approx \sqrt{\frac{2\pi}{n\epsilon''(\mathbf{w}^*)}} e^{-n\beta\epsilon(\mathbf{w}^*)}, \quad (67)$$

1337 where  $\mathbf{w}^*$  satisfies the saddle point equations:

$$\epsilon(\mathbf{w}^*) := \frac{\partial}{\partial \mathbf{w}} \epsilon(\mathbf{w})|_{\mathbf{w}=\mathbf{w}^*} = 0 \quad (68)$$

$$\epsilon(\mathbf{w}^*) := \frac{\partial^2}{\partial^2 \mathbf{w}} \epsilon(\mathbf{w})|_{\mathbf{w}=\mathbf{w}^*} > 0. \quad (69)$$

1338 To apply the SPA rigorously, we expect that  $\epsilon(\mathbf{w})$  decays exponentially, whereas we will be  
 1339 working with Heavy-Tailed (HT)/Power-Law (PL) distributions. At first glance, this might seem  
 1340 problematic. These distributions, however, can also be modelled as Truncated Power Laws (TPL)  
 1341 because of the finite-size effects. This suggests we can apply the SPA formally, assuming (but not  
 1342 checking) that the finite-size effects will satisfy the requirements we need. The SPA will play a  
 1343 central role when we evaluate the matrix-generalized form of our model and the HCIZ integrals  
 1344 that arise.

**When the limit fails: atypical behavior.** When the Thermodynamic limit fails to exist, the Free Energy will contain additional, non-extensive terms, i.e.,  $F = n\bar{F}_{ex} + n^{1+x}\bar{F}_{non-ex}$ ,  $x > 0$ .<sup>26</sup> In this case, the AA may fail, the SPA may not apply, and the system may fail to be Self-Averaging. This causes the system to behave in an atypical way, possibly converging to a meta-stable and/or glassy phase. Indeed, when the weights  $\mathbf{w}$  are *atypical*, they may describe the training data well, but they would fail to describe the test data well; in this sense, we say the model is *overfit* to the training data. We will not explicitly consider a model that is non-extensive; however, we will present empirical results where we suspect the model is overfit (in Section 6.4) and, additionally, where we observe glassy behavior, which we refer to as a *Hysteresis Effect* (in Section 6.5).

**Remark.** Additional results are provided in Appendix A.2. In particular, in Appendix A.2.1, we use the ideas from this section to derive the full non-linear vector form of the Annealed Hamiltonian  $H^{an}(\mathbf{w})$  (Eqn. 44) for the Linear Perceptron, in the AA, and then we express it at High-T, i.e.,  $H_{hT}^{an}(R) = \epsilon(R)$  (Eqn. 49). Then, in Appendix A.2.1, we derive the matrix generalization  $H^{an}(R) \rightarrow H^{an}(\mathbf{R})$  of these quantities. This derivation is necessary to obtain the normalization on the weight matrices  $\mathbf{W}$  necessary for the final results in Section 5.

#### 4.2.8 HCIZ Integrals

To generalize the Linear ST Perceptron (in the AA, and at high-T) from Perceptron vectors to MLP matrices, we need to generalize the thermal average over the  $m$ -dimensional Perceptron weight vectors ( $\mathbf{w} = \mathbf{s}$ ) to an integral over NN Student  $N \times M$  weight matrices ( $\mathbf{W} = \mathbf{S}$ ). The resulting Partition Function and Free Energy will be expressed with what is called an HCIZ integral. (See Tanaka [?, ?], Gallucio et al. [?] and/or Parisi [?], and also Eqn. 129 in Section 5.) Analogously to Eqn. ??, we will define a *Layer Quality Generating Function*,  $\beta\Gamma_{\mathcal{Q}^2}^{IZ}$ , for the Layer Quality (squared). This will take the form of an HCIZ integral,

$$\beta\Gamma_{\mathcal{Q}^2}^{IZ} := \lim_{N \gg 1} \ln \underbrace{\int d\mu(\mathbf{A}) \exp[N\beta \text{Tr}[\mathbf{A}_2 \mathbf{X}_2]]}_{\text{HCIZ Integral}} \approx N\beta \text{Tr}[\mathcal{G}_A(\lambda)], \quad (70)$$

that is evaluated in the “large-N limit” (here, meaning as  $N \gg 1$ , but not at the limit  $N = \infty$ ). Here,  $\mathbf{A}$  and  $\mathbf{X}$  are  $N \times N$  Hermetian matrices,  $d\mu(\mathbf{A})$  is a measure over all random (Orthogonal) matrices (see Eqn. 243), and  $\mathcal{G}_A()$  is a (Norm)? Generating Function (different from  $\Gamma$ , above), defined in Eqn. 132 in Section 5. In applying this, we will actually write  $\text{Tr}[\mathbf{A}\mathbf{X}] = \frac{1}{N} \text{Tr}[\mathbf{A}\mathbf{W}\mathbf{W}^\top] = \frac{1}{N} \text{Tr}[\mathbf{W}^\top \mathbf{A}\mathbf{W}]$ , where  $\mathbf{W}$  is an  $N \times M$  layer weight matrix, and  $\mathbf{A} = \mathbf{A}_2 := \frac{1}{N} \mathbf{S}\mathbf{S}^\top$  is a layer Correlation matrix, and, here,  $\mathbf{X} = \mathbf{X}_2 = \frac{1}{N} \mathbf{W}\mathbf{W}^\top$ . Moreover, to evaluate this, we will need to restrict  $\mathbf{A}$  (and  $\mathbf{X}$ ) to the lower-rank Effective Correlation Space (ECS), i.e.,  $\mathbf{A} \rightarrow \tilde{\mathbf{A}}$ .

Notice this looks similar to Saddle Point Approximation (SPA), but where the more complicated function  $\mathbb{G}_{\mathbf{A}}()$  now appears. Also,  $\mathbb{G}_{\mathbf{A}}()$  here depends on only the limiting form of the ESD of  $\mathbf{A}$ ,  $\rho_A^\infty(\lambda)$ , and depends on the  $M$  normalized eigenvalues  $\lambda_i$  of  $\mathbf{X}$ .

To evaluate this, we will form the large-N limit, using a result from Tanaka [?, ?], but extended slightly (in Appendix A.6) to include the inverse-Temperature  $\beta$  explicitly. We can write

$$\beta\Gamma_{\mathcal{Q}^2, N \gg 1}^{IZ} := \lim_{N \gg 1} \beta\Gamma_{\mathcal{Q}^2}^{IZ}, \quad (71)$$

<sup>26</sup>When dealing with matrix integrals (below),  $F \sim NMF_0 + \dots$ , when there are  $N \times M$  degrees of freedom [?], but we will only be concerned with the large-N limit.

1380 with the final result

$$\beta \mathbf{T}_{\mathcal{Q}^2, N \gg 1}^{IZ} := N\beta \sum_{i=1}^{\tilde{M}} \int_{\lambda_{min}^{ECS}}^{\tilde{\lambda}_i} dz R(z), \quad (72)$$

1381 where  $R(z)$  is a complex function, the R-transform of the ESD  $\rho(\lambda)$  of the Teacher, and  $\tilde{\lambda}_i$  are the  
 1382 eigenvalues of Teacher Correlation Matrix  $\tilde{\mathbf{X}}$ , restricted to the ECS. For more details, see Section 5  
 1383 and Appendices A.4 and A.6.

1384 **Branch Cuts and Phase Behavior.** Free energies  $(F, \Gamma)$  often exhibit *branch cuts* when  
 1385 expressed as analytic functions of complex parameters (i.e., temperature, coupling constants, or  
 1386 eigenvalue cutoffs), and arise from singularities in the underlying integral representations of the  
 1387 partition function  $Z$ . When a branch cut occurs, this demarcates non-analytic behavior and this  
 1388 indicates the onset of a *phase transition* where macroscopic observables such as correlation lengths  
 1389 and/or variance-like quantities may diverge or change character abruptly.

1390 In the context of our HCIZ-based construction, integrating a complex function like the  $R(z)$ -  
 1391 transform of a Heavy-Tailed ESD can produce precisely this phenomenon. For example, if  $R(z)$   
 1392 has a square-root term, i.e.,  $(\sqrt{z-c})$ , then it will have a branch cut at  $z=c$ , and the Generating  
 1393 Function (i.e., Effective Free Energy)  $\beta \mathbf{T}_{\mathcal{Q}^2, N \gg 1}^{IZ}$  will be non-analytic and we must choose the  
 1394 appropriate, physically meaningful branch, i.e.,  $(z > c)$ , corresponding to the ECS. We argue that  
 1395 this cut signifies a *phase boundary*—an abrupt change in the system’s correlation structure and  
 1396 corresponds to an emerging singularity in the Layer Quality. Even though we perform only a single  
 1397 exact RG step, (rather than fully iterating a renormalization flow), the appearance of a branch cut  
 1398 in  $\beta \mathbf{T}_{\mathcal{Q}^2, N \gg 1}^{IZ}$  will encode non-trivial *phase-like* behavior in the SETOL Heavy-Tailed matrix model.

### 1399 4.3 Student -Teacher Model

1400 In this subsection, we present a unified view of the *Student-Teacher* (ST) model from both a  
 1401 practical and a theoretical perspective. From the practical (*operational*) side, imagine a real-world  
 1402 practitioner training a large-scale neural network (NN); call this NN the Teacher. We wish to assess  
 1403 the Teacher’s *true* accuracy on ground-truth labels in consistently reproducing its own (possibly  
 1404 imperfect) outputs, one can train another network—the Student—to mimic the Teacher’s predictions.  
 1405 By comparing Student and Teacher outputs, one can approximate the Teacher’s generalization  
 1406 performance even without an explicit hold-out set. Notably, if the Teacher perfectly interpolates  
 1407 its training data, the Student’s error directly estimates the Teacher’s *true* Generalization Accuracy.  
 1408 Otherwise, it captures the Teacher’s *Precision* in reproducing its own noisy or biased labels.

1409 From the theoretical side, we seek a succinct, analytic formulation of the ST Average General-  
 1410 ization Error, denoted  $\tilde{\mathcal{E}}_{gen}^{ST}$ . We work in the Annealed Approximation ( $\tilde{\mathcal{A}}$ )—a simplification often  
 1411 valid when networks are sufficiently large and can nearly interpolate their training data. Under this  
 1412 approximation, one obtains closed-form expressions for the Student-Teacher *Overlap* ( $R$ ) and thus  
 1413 for the Teacher’s overall error or accuracy. These results lay the foundation for our *Semi-Empirical*  
 1414 approach, in which we supplement this theoretical form with empirical measurements (e.g., from  
 1415 the trained weight matrices) to account for real-world correlations in the data and the model’s  
 1416 internal structure.

1417 **4.3.1 Operational Setup** In subsection 4.3.1 we explain how to set up the classic Student-Teacher  
 1418 model in an *operational* manner. In particular emphasize the difference between *true accuracy*  
 1419 (vs. ground-truth labels) and Precision (vs. the Teacher’s own labels). We also discuss  
 1420 the difference between our Quality metrics and the *GeneralizationGap*.

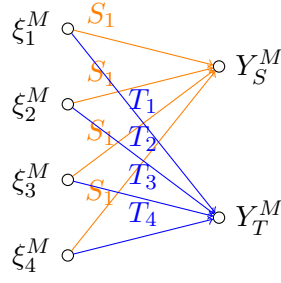


Figure 7: Pictorial representation Student and Teacher Perceptrons.

### 4.3.2 Theoretical Student-Teacher Average Generalization Error ( $\bar{\mathcal{E}}_{gen}^{ST}$ )

In subsection 4.3.2 we outline how to derive  $\bar{\mathcal{E}}_{gen}^{ST}$  using the AA. We introduce the key expressions that will serve as the starting point for our extended Semi-Empirical theory in subsequent subsubsection.

#### 4.3.1 Operational Setup

Here, describe the basic setup of the classic Student-Teacher model, taking an operational view from the perspective of a practitioner training real-world Student and Teacher models. Specifically, we present the Annealed Approximation (AA) in a practical light, and use it explain the difference between computing the *Empirical Generalization Error*,  $\bar{\mathcal{E}}_{gen}^{emp}$ , for the *True Accuracy* and for the *Precision* of a Teacher model.

**Test Error of the Teacher** We start by describing how to obtain a simple formal expression for the empirical test errors of the Teacher, first for the True Accuracy.

Let us say we have a model, called Teacher ( $T$ ), which maps some *actual* (i.e., correlated) data  $\mathbf{x}_\mu \in \mathbf{D}$  to some known or *true* labels ( $\mathbf{y}_\mu^{true}$ ) (where,  $\mathbf{y}_\mu^{true}$  is, say, an  $N$ -dimensional vector of binary labels). We might say that  $\mathbf{y}_\mu^{true}$  represents the *Ground Truth* for the problem. Operationally, we train the Teacher  $T$  to reproduce or at least approximate the true labels  $\mathbf{y}_\mu^{true}$ .

$$T : \mathbf{x}_\mu \rightarrow \mathbf{y}_\mu^T \approx \mathbf{y}_\mu^{true}. \quad (73)$$

If  $T$  reproduces the true labels exactly, we might say that the Teacher has been trained to *Interpolation*, and, therefore,  $\mathbf{y}_\mu^T = \mathbf{y}_\mu^{true}$ . Indeed, most models today are trained to *Interpolation*, and we don't need to necessarily worry about the difference between the true and the predicted Teacher labels. Formally, however, and to understand better the  $\hat{\Delta}$ , it is beneficial to discuss the implication of this distinction.

Following Eqn. 1, let's say the Teacher outputs are specified by an Energy function  $E_{NN}^T$

$$\mathbf{y}_\mu^T = E_{NN}(T, \mathbf{x}_\mu) \quad (74)$$

<sup>27</sup> so that we may write the *Empirical Average Training Error*  $\bar{\mathcal{E}}_{train}^{emp}$  as

$$\bar{\mathcal{E}}_{train}^{emp} := \frac{1}{N^{train}} \sum_{\mu=1}^{N^{train}} \mathcal{L}[\mathbf{y}_\mu^{train}, E_{NN}(T, \mathbf{x}_\mu^{train})]. \quad (75)$$

<sup>27</sup>Do not confuse the Energy/output function  $E_{NN}$  with the energies  $\Delta \mathbf{E}$  defined below to represent the ST error function(s). We refer to outputs of  $E_{NN}(\mathbf{t}, \mathbf{x}_\mu)$ , when applied to a data point  $\mathbf{x}_\mu$ , as energies because they are effectively un-normalized probabilities for the class outputs (for labels  $y_\mu = 1$  or  $-1$ ).



1443 Ideally, we seek the *True Average Generalization Error* of the Teacher, denoted  $\bar{\mathcal{E}}_{gen}^T$ . Of course, this  
 1444 is unknowable, but in practice, we estimate  $\bar{\mathcal{E}}_{gen}^T$  by measuring the error of the Teacher predictions  
 1445 on some test (or hold-out) set  $(\mathbf{x}_\mu^{test}, y_\mu^{test})$ . We call this the *Empirical Average Generalization*  
 1446 *Error*,  $\bar{\mathcal{E}}_{gen}^{emp}$ , and write

$$\bar{\mathcal{E}}_{gen}^T \approx \bar{\mathcal{E}}_{gen}^{emp} := \frac{1}{N^{test}} \sum_{\mu=1}^{N^{test}} \mathcal{L}[y_\mu^{test}, E_{NN}(T, \mathbf{x}_\mu^{test})]. \quad (76)$$

1447 To measure the error, the loss function  $\mathcal{L}$  may be a L1 ( $\ell_1$ ) or L2 ( $\ell_2$ ) loss; whereas for training a  
 1448 model, it is usually something like a cross-entropy loss, but this detail does matter later.

1449 If we don't have a hold-out set, however, we can still estimate  $\bar{\mathcal{E}}_{gen}^T$  using the Student-Teacher  
 1450 approach.

1451 **Estimating the Teacher Error with Students: Accuracy vs. Precision** Imagine training  
 1452 a Student ( $S$ ) model (with a similar architecture as  $T$ , and acting on the same dataset  $\mathbf{x}_\mu \in \mathbf{D}$ ),  
 1453 which tries to reproduce the Teacher predictions:

$$S : \mathbf{x}_\mu \rightarrow y_\mu^S \approx y_\mu^T, \quad (77)$$

1454 and assume the Student outputs  $y_\mu^S$  are given by the Energy output function  $E_{NN}^S$

$$y_\mu^S = E_{NN}(S, \mathbf{x}_\mu), \quad (78)$$

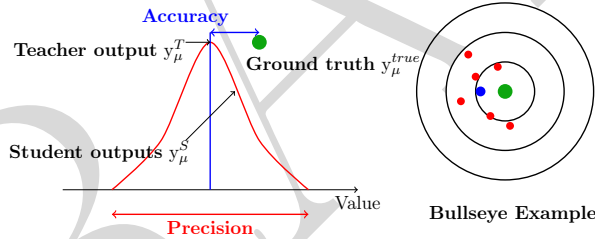


Figure 8: Precision vs. Accuracy

1455 If the Teacher is trained to Interpolation, then the difference between the Student and the  
 1456 Teacher labels estimates the true error, i.e.,  $\|y_\mu^S - y_\mu^T\| = \|y_\mu^S - y_\mu^{true}\|$ , and this error is associated  
 1457 with the Training Accuracy of the model in predicting the Ground Truth. But if the Teacher makes  
 1458 some errors, then  $\|y_\mu^S - y_\mu^T\|$  is now estimating the Precision of the model. These two situations are  
 1459 depicted in Figure 8.

1460 The Student-Teacher model also explains why NNs can generalize even when trained  
 1461 to Interpolation on noisy data (which has been a source of confusion [?]). In this model, the  
 1462 Generalization Error  $\bar{\mathcal{E}}_{gen}^T$  is a simple function of the overlap  $R$  between the Teacher  $T$  and the  
 1463 Students  $S$ , i.e.,  $\bar{\mathcal{E}}_{gen}^T \sim \langle 1 - \epsilon(R) \rangle_s^\beta$ . So even if the Teacher  $T$  is trained on noisy data, as long as  
 1464 there are Students  $S$  with significant overlap  $R$  with the Teacher, the Teacher Generalization Error  
 1465  $\bar{\mathcal{E}}_{gen}^T$  can be considerably small. For more details, see [?]

1466 **Learning the Student** Moving forward, we will always assume the Teacher is trained to  
 1467 Interpolation because this actually corresponds to the Annealed Approximation, whereas if the  
 1468 Teacher makes errors, we may need to consider Quenched averages, explained below.

Imagine now that in order to estimate the empirical Average Generalization Error,  $\bar{\mathcal{E}}_{gen}^{emp}$ , by training a very large number of Students, and computing the average ST error on some test set. Let us break the data set into training  $\mathbf{x}_\mu^{train}$  and test  $\mathbf{x}_\mu^{test}$  examples, train models on the training data (that is, find the optimal model weights), and evaluate the  $S$  and  $T$  models on the test data.

The Student learning task can be written as in Eqn. (2) as the following optimization problem over the training data:

$$\operatorname{argmin}_{\{S'\}} \sum_{\mu=1}^{N^{train}} \mathcal{L}[E_{NN}(S', \mathbf{x}_\mu^{train}), E_{NN}(T, \mathbf{x}_\mu^{train})], \quad (79)$$

If the Teacher is trained to Interpolation, then the optimization problem in Eqn. ?? is training a Student to reproduce the Ground Truth labels, so that  $y_\mu^S \sim y_\mu^{true}$  for both the training and test sets.

$$\operatorname{argmin}_{\{S'\}} \sum_{\mu=1}^{N^{train}} \mathcal{L}[E_{NN}(S', \mathbf{x}_\mu^{train}), y_\mu^{true}], \quad (80)$$

But if not, then the Student is reproducing the possibly incorrect Teacher labels, and, importantly, the Student  $S$  now depends explicitly on how the Teacher was trained. That is, we should denote that the learned Student explicitly depending on  $T$ , i.e.  $S \rightarrow S[T]$ . This will be important below.

**The Average Generalization Error** In either case, however, we may still estimate the Empirical Average Generalization Error by replacing the test predictions in Eqn. 76 with the student predictions  $y_\mu^{test} \rightarrow y_\mu^S$ , and then averaging directly over the test data  $\mathbf{x}_\mu^{test}$  for all possible or available test examples.

If we have a very large number of suitable Students (say, drawn from some random distribution), then we can try to estimate the Average Generalization Error of the Teacher, i.e.,  $\mathcal{E}_{gen}^T \approx \bar{\mathcal{E}}_{gen}^{emp}$ .  $\bar{\mathcal{E}}_{gen}^{emp}$  is given by an average loss, the average is over all possible Students  $N_S$ , and then over all  $N^{test}$  test data points  $\mathbf{x}_\mu^{test} \in \mathbf{D}$

$$\begin{aligned} \bar{\mathcal{E}}_{gen}^{emp} &= \frac{1}{N^{test}} \sum_{\mu=1}^{N^{test}} \frac{1}{N_S} \sum_S \mathcal{L}[E_{NN}(S, \mathbf{x}_\mu^{test}), E_{NN}(T, \mathbf{x}_\mu^{test})] \\ &= \frac{1}{N_S} \sum_S \frac{1}{N^{test}} \sum_{\mu=1}^{N^{test}} \mathcal{L}[E_{NN}(S, \mathbf{x}_\mu^{test}), E_{NN}(T, \mathbf{x}_\mu^{test})], \end{aligned} \quad (81)$$

where (ideally)  $N^{test}$  is extremely large.

In Bra-Ket notation, we may also write

$$\begin{aligned} \bar{\mathcal{E}}_{gen}^{emp} &= \langle \langle \Delta \mathbf{E}_{\mathcal{L}}(S, T, \mathbf{x}) \rangle_S \rangle_{\mathbf{x}^{test}} \\ &= \langle \langle \Delta \mathbf{E}_{\mathcal{L}}(S, T, \mathbf{x}) \rangle_{\mathbf{x}^{test}} \rangle_S \end{aligned} \quad (82)$$

where  $\Delta \mathbf{E}_{\mathcal{L}}(S, T, \mathbf{x}) := \mathcal{L}[E_{NN}(S, \mathbf{x}_\mu^{test}), E_{NN}(T, \mathbf{x}_\mu^{test})]$ . For the empirical estimate, it does not matter what order we take the sums in, but we are not estimating the the True Average Generalization Error of the Teacher,  $\bar{\mathcal{E}}_{gen}^T$ , unless  $T$  is trained to Interpolation. For the theoretical estimate, however, the order can be important, and this also depends on if  $T$  is trained to Interpolation or not. <sup>28</sup>

<sup>28</sup>This approach can be likened to the Bootstrap method [?] used for error estimation. However, the Bootstrap method predominantly emphasizes variations in the input data  $\mathbf{x}^n \in \mathbf{D}$ , while in this context, we are essentially bootstrapping over the students  $S$ .

**Annealed vs. Quenched Averages** [THIS NEEDS CLEANED UP] Recall that in the **StatMech** approach to computing errors, we do not break the data into training and test, but, instead, to obtain the Average Generalization Error,  $\bar{\mathcal{E}}_{gen}$ , use the Generating Function approach. In doing this, we need to compute both the Thermal Average over the model weights  $(S, T)$ , and also take the data average over the entire available model data set  $\mathcal{D}$ . And the order can matter.

In the case where the Teacher is trained to Interpolation, may may train the Student independently of *when* the Teacher. But if Teacher is *not* trained to Interpolation, then formally we must train the Teacher first to obtain target predictions for the Student. That is, the Student formally depends on the Teacher,  $S[T]$ . The empirical errors in  $T$  would then formally depend on the specific instantiation of the data  $\xi^n \in \mathcal{D}$ , and therefore, conceptually imagine that we must first average over the data before averaging over the weights. Training to Interpolation corresponds conceptually to working with a model in the Annealed Approximation, whereas not corresponds to Quenched case.

Practically, when the Teacher is not trained to Interpolation, we may need to resample the training data and training an ensemble of models to compensate for anomalies in training data (bad labels, noise, etc.) that may cause the underlying model to overfit to the training data. Theoretically, within **SMOG**, this is equivalent to *quenching* the system to the data (a term analogous to quickly cooling a physical system, freezing in any defects). In contrast, when one *anneals* a physical system, one heats up and cools it down slowly, and repeatedly, thereby removing any defects (of data anomalies for NNs, or material defects in a physical system).

In **StatMech**, one can perform a so-called quenched average using a replica calculation, effectively removing the dependence on test and/or training data from the final estimate for  $\bar{\mathcal{E}}_{gen}^{emp}$ . However, the theoretical quenched result may differ significantly from the annealed case when the underlying model is unrealizable [?]. This may occur when the training data is very noisy and/or the model architecture is such that it can not correctly predict all the training labels. In such cases, the model will always have some finite, non-zero Average Training Error,  $\bar{\mathcal{E}}_{train} > 0$ , even in the large- $N$  limit of infinite data. In such a case, this indicates a highly complex error landscape with many local minima separated by extremely high barriers, and a slowing down of the dynamics.<sup>29</sup>

While it is commonplace to train ensembles and/or use cross-validation when training small models (as the above discussion assumes), this could be extremely expensive and impractical in modern ML, e.g., for very large models like LLMs [?]. For such massive NNs, one needs a theory that can detect anomalies in training directly from observations during and/or after training. This is a hallmark of the **SETOL** approach, and it distinguishes **SETOL** from the classic **StatMech** approach.

be used to estimate the Average Generalization Accuracy (and not the Precision), and which we will refer to more generally as a layer and/or model Quality.

**Generalization Gap vs. Model Quality** We should distinguish between what is typically done in the **SLT** literature versus the **StatMech** approaches. In **SLT**, one is typically interested in modeling the *Generalization Gap*. The Generalization Gap quantifies the difference between a models performance on training data versus unseen test data:

$$\mathcal{E}_{gap}^{emp} := \mathcal{E}_{train}[\mathbf{x}^{train}] - \mathcal{E}_{gen}[\mathbf{x}^{test}] \quad (83)$$

In contrast, in **StatMech** approaches, one considers the Model Generalization Accuracy directly, which is sometimes called the Model Quality in the **SLT** literature. Model quality is an indication of the models accuracy, precision, recall, or any other relevant metric based on the task at hand.

<sup>29</sup>In modern ML parlance, one might say the model can not be evaluated at interpolation, although in practice such a model might have a zero empirical Training Error since it may overfit the specific training data.

While related, in developing an analytic theory, the Generalization Gap and the Model Quality (or Model Generalization Accuracy) require conceptually different approaches. This is because the Generalization Gap depends on a specific realization of the training data, whereas our Model Generalization Accuracy will be formulated on a random training data set (and then corrected later with empirical data). In this sense, any theory of the Generalization Gap requires a formalism where the predicted model error is Quenched to the training data, which is not what we want. In contrast, the Model Generalization Accuracy will be formulated using the Annealed Approximation (AA), and is therefore both conceptually and mathematically simpler.

[\[Comment on our paper with YQ\]](#)

### 4.3.2 Theoretical Student-Teacher Average Generalization Error ( $\bar{\mathcal{E}}_{gen}^{ST}$ ),

Here, we seek a simple, formal expression for the Student-Teacher Average Generalization Error,  $\bar{\mathcal{E}}_{gen}^{ST}$ , that can be used as the starting point for our extended Semi-Empirical theory.

**The Data Model.** To develop a Semi-Empirical theory of the Teacher Generalization Error,  $\mathcal{E}_{gen}^T$ , instead of training and evaluating a NN model using real data ( $\mathbf{x}$ ), we seek a simple, analytical expression with parameters that can be fit to empirical measurements. So in addition to using a model for our NN, we must specify a model for the data. In a real NN, the data  $\mathbf{x}$  is correlated, and, in fact, very strongly correlated; and this is reflected in the layer weight matrices. However, to be tractable, our starting theoretical expressions use uncorrelated (i.i.d) data. Formally, we must replace the correlated data with some uncorrelated, random model of the data, i.e.,  $\mathbf{x} \rightarrow \boldsymbol{\xi}$ . As described in Figure 6, our Data Model is a standard Gaussian  $N(0,1)$  model for the input data

$$\mathbf{x}_\mu \rightarrow \boldsymbol{\xi}_\mu, \quad \boldsymbol{\xi}_\mu \in N(0,1), \quad (84)$$

where  $N(0,1)$  denotes a Gaussian distribution with zero mean and unit variance,  $\boldsymbol{\xi}_\mu$  is normalized such that  $\|\boldsymbol{\xi}_\mu\|_F^2 = \frac{1}{M}$  for all  $N$  data vectors.

We make this distinction between Actual and Model Data to emphasize that, later, we will use our so-called Semi-Empirical procedure to account for the real correlations in the actual data phenomenologically by taking some analytical parameter of the theory and fitting it to the real world observations, here, on the ESD of the NN weight matrices.

**The ST Error Model and the Annealed Potential  $\epsilon(\mathbf{s}, \mathbf{t})$ .** We now model Teacher error  $\bar{\mathcal{E}}_{gen}^T$  with the *Average ST Generalization Error*  $\bar{\mathcal{E}}_{gen}^{ST}$ , which is obtained by *first* computing the ST error function  $\Delta \mathbf{E}_{\mathcal{L}}(\mathbf{s}, \mathbf{t}, \boldsymbol{\xi})$  over the set of *all* possible  $N$  input examples  $\boldsymbol{\xi}$ . Define the data-dependent ST test error function –or Energy difference– as

$$\Delta \mathbf{E}_{\mathcal{L}}(\mathbf{s}, \mathbf{t}, \boldsymbol{\xi}) := \sum_{\mu=1}^N \mathcal{L}[E_{NN}(\mathbf{s}, \boldsymbol{\xi}_\mu), E_{NN}(\mathbf{t}, \boldsymbol{\xi}_\mu)]. \quad (85)$$

where  $\mathcal{L}(\mathbf{s}, \mathbf{t}, \boldsymbol{\xi})$  is simply the  $\ell_2$  loss. This measures the error between the Student and the Teacher; it is zero when their predictions are identical, ( $\mathbf{y}_\mu^S = \mathbf{y}_\mu^T$ ) when ( $\boldsymbol{\xi} = \boldsymbol{\xi}_\mu$ ), and is nonzero otherwise.

We aim to derive a simple expression for the Average ST Generalization Error,  $\bar{\mathcal{E}}_{gen}^{ST}$ , and to do this, we define the Annealed Error Potential for the data-averaged ST Generalization Error  $\epsilon(\mathbf{s}, \mathbf{t})$ , as in Eqn. 30, as:

$$\epsilon(\mathbf{s}, \mathbf{t}) = \langle \Delta \mathbf{E}_{\mathcal{L}}(\mathbf{s}, \mathbf{t}, \boldsymbol{\xi}) \rangle_{\boldsymbol{\xi}^N} := \frac{1}{N} \int d\mu(\boldsymbol{\xi}^N) \Delta \mathbf{E}_{\mathcal{L}}(\mathbf{s}, \mathbf{t}, \boldsymbol{\xi}) \quad (86)$$

1575 The measure  $d\mu(\boldsymbol{\xi}^N)$  will end up being a Gaussian measure over  $N$  samples (see Appendix A.2),  
 1576 and the intent is to evaluate it in the large- $N$  limit, thereby sampling all possible inputs in the  
 1577 model space,  $\boldsymbol{\xi} \in \mathcal{D}$ .

1578 As in Eqn. ?? (Section 4.2), by applying the AA, we can rewrite the Average ST Generalization  
 1579 Error,  $\bar{\mathcal{E}}_{gen}^{ST}$ : first, a simple average over all the possible inputs  $\boldsymbol{\xi}$ ; and, second, then as a Thermal  
 1580 average over all Students  $S$ , in the AA, and at high- $T$

$$\bar{\mathcal{E}}_{gen}^{ST} := \langle \epsilon(\mathbf{s}, \mathbf{t}) \rangle_{\mathbf{s}}^{\beta}. \quad (87)$$

1581 (Recall that in this regime  $\bar{\mathcal{E}}_{gen}^{ST} = \bar{\mathcal{E}}_{gen}^{an, hT}$ .)

1582 In the classic **StatMech** approach, the average  $\langle \dots \rangle_{\mathbf{s}}^{\beta}$  is a Thermal Average in the canonical  
 1583 ensemble with  $\beta$  fixed, as explained in Section 4.2. Here, we will do something similar, as the  
 1584 Student average  $\langle \dots \rangle_{\mathbf{s}}^{\beta}$  will be computed from the associated generating function  $\beta \mathbf{T}_{\mathcal{Q}^2}^{IZ}$  for the  
 1585 matrix-generalized case (i.e., an HCIZ integral defined over all students, and in the large- $N$  limit).)

1586 Recall that above, the empirical estimate for  $\bar{\mathcal{E}}_{gen}^{emp}$  depended on a specific instantiation of the  
 1587 model for the training data  $\mathbf{x}_{\mu}^{train}$ , i.e.  $\bar{\mathcal{E}}_{gen}^{ST}$  is Quenched to the training data. For that reason, for  
 1588 the final result, we needed to take a second, quenched average over all possible data sets. Here,  
 1589 we do not need to consider this and always work in the Annealed Approximation (AA). This is  
 1590 because we incorporate the specific effects of the real-world training data ( $\mathbf{x}^n$ ) after we derive our  
 1591 formal expressions by fitting the model parameters to empirical data. The final expression for  $\bar{\mathcal{E}}_{gen}^{ST}$ ,  
 1592 derived below, will be generalized to  $\bar{\mathcal{E}}_{gen}^{NN}$ , matrix-generalization of the classic **StatMech** formula  
 1593 for the Linear Perceptron, in the Annealed and High- $T$  approximations. (see Appendix A.2).

1594 **The Annealed Potential as a function of the overlap ( $\epsilon(R)$ ).** We want an expression for  
 1595 the data average of the ST test error, from Eqn. 86, generalized from Perceptron vectors to NN  
 1596 layer weight matrices. For the Perceptron, one obtains different expressions for the ST error  
 1597 function, depending on the type of activation function  $h(x)$  in Eqn. ??; The simplest are the  
 1598 Linear and Boolean Perceptrons, and for both (and with  $\ell_2$  loss),  $\epsilon(\mathbf{s}, \mathbf{t})$  is simply a function of the  
 1599 ST overlap  $R$  [?]. This gives  $\epsilon(\mathbf{s}, \mathbf{t}) \rightarrow \epsilon(R)$ , where

$$R = \frac{1}{N} \mathbf{s}^{\top} \mathbf{t} = \frac{1}{N} \sum_{i=1}^M s_i t_i, \quad (88)$$

1600 which is simply the dot product between the  $M$ -dimensional Student  $\mathbf{s}$  and Teacher  $\mathbf{t}$  weight  
 1601 vectors, and normalized by the number of training instances  $N$ . For a Linear Perceptron [?],<sup>30</sup>  
 1602 with activation function  $h(x) = x$ , the error function is

$$\epsilon(R) = 1 - R. \quad (89)$$

1603 **Derivation of the ST error ( $\epsilon(R)$ ) for the Linear Perceptron.** To derive Eqn. 89, define  
 1604 the data-dependent ST error (Eqn. 85) in terms of an  $\ell_2$  loss function

<sup>30</sup>In the classic approach for the ST model, the theory examined different expressions  $\epsilon(R)$ . For example, one can consider the Boolean Perceptron [?, ?], with activation function  $h(x) = \text{sgn}(x)$ , i.e., the Heaviside step function. Then, the error is  $\epsilon(R) = 1 - \frac{1}{\pi} \arccos(R)$ . In both cases, perfect learning occurs when  $R = 1$  [?].

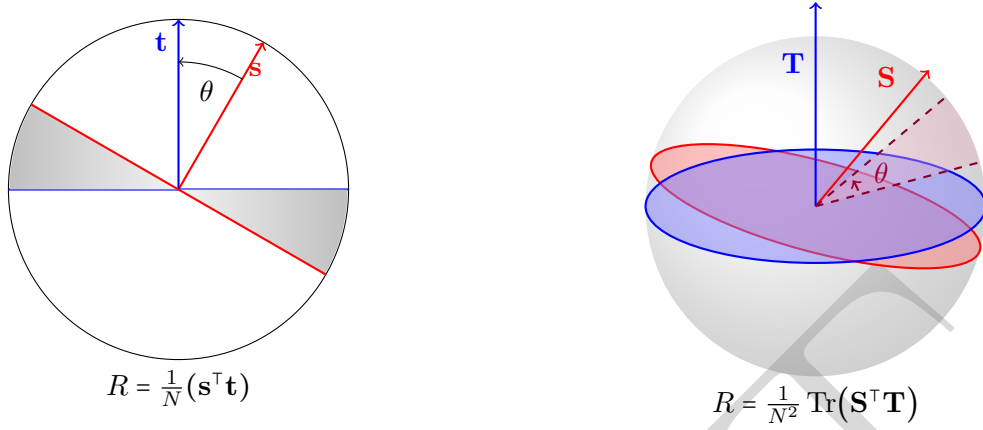


Figure 9: Comparison of 2D and 3D representations of the vector and matrix Student-Teacher overlap  $R$ . **Left:**  $R = \frac{1}{N} \mathbf{s}^\top \mathbf{t}$ . **Right:**  $R = \frac{1}{N^2} \text{Tr}(\mathbf{S}^\top \mathbf{T})$  with conic sections on the sphere (red  $\mathbf{S}$ , blue  $\mathbf{T}$ ), plus a purple wedge for the angle.

$$\begin{aligned}
\Delta \mathbf{E}_{\ell_2}(\mathbf{s}, \mathbf{t}, \boldsymbol{\xi}^n) &= \frac{1}{2} (\mathbf{y}^S - \mathbf{y}^T)^\top (\mathbf{y}^S - \mathbf{y}^T) \\
&= \frac{1}{2} [(\mathbf{y}^S)^\top \mathbf{y}^S - 2(\mathbf{y}^S)^\top \mathbf{y}^T + (\mathbf{y}^T)^\top \mathbf{y}^T] \\
&= N - (\mathbf{y}^S)^\top (\mathbf{y}^T) \\
&= N(1 - \eta(\mathbf{s}, \mathbf{t})),
\end{aligned} \tag{90}$$

where we call  $\eta(\mathbf{s}, \mathbf{t}, \boldsymbol{\xi}) := \frac{1}{N} \mathbf{y}_S^\top \mathbf{y}_T$  the *data-dependent Self-Overlap*. The expression  $\eta(\mathbf{s}, \mathbf{t}, \boldsymbol{\xi})$  is analogous to the ST overlap  $R$ , but before the data has been integrated out. The Self-Overlap  $\eta(\mathbf{s}, \mathbf{t}, \boldsymbol{\xi})$  will appear later in Eqn. ?? (in Section 5), when formulating the matrix-generalized overlap operator  $\mathbf{R}$ .

Using the  $E_{NN}$  Energy generating or output function (Eqn. ??), we can replace the label vectors  $\mathbf{y}^T, \mathbf{y}^S$  as

$$\mathbf{y}^S = \mathbf{s}^\top \boldsymbol{\xi}, \quad \mathbf{y}^T = \mathbf{t}^\top \boldsymbol{\xi}, \tag{91}$$

which gives

$$\eta(\mathbf{s}, \mathbf{t}, \boldsymbol{\xi}) = \frac{1}{N} (\mathbf{s}^\top \boldsymbol{\xi})^\top (\mathbf{t}^\top \boldsymbol{\xi}) = \frac{1}{N} \boldsymbol{\xi}^\top (\mathbf{s}^\top \mathbf{t}) \boldsymbol{\xi}. \tag{92}$$

or, more simply,

$$\eta(\mathbf{s}, \mathbf{t}) = \langle \eta(\mathbf{s}, \mathbf{t}, \boldsymbol{\xi}) \rangle_{\boldsymbol{\xi}^N} = \frac{1}{N} (\mathbf{s}^\top \mathbf{t}) \tag{93}$$

We want to evaluate this as an Annealed Error Potential  $\epsilon(R)$  for the data-averaged ST test error, as in Eqn. 86. To do this, we need to compute the average or Expected Value over all  $N$



possible input data vectors  $\xi$  (i.e.,  $d\mu(\xi^N) = \mathcal{D}\xi^N P(\xi^N)$ ).

$$\langle \Delta \mathbf{E}_{\ell_2}(\mathbf{s}, \mathbf{t}, \xi^n) \rangle_{\xi^N} = \int d\mu(\xi^N) (1 - \eta(\xi)) \quad (94)$$

$$\begin{aligned} &= \int d\mu(\xi^N) (1 - \xi^\top \frac{1}{N} \mathbf{s}^\top \mathbf{t} \xi) \\ &= \int d\mu(\xi^N) - \int d\mu(\xi^N) \xi^\top \frac{1}{N} \mathbf{s}^\top \mathbf{t} \xi \\ &= 1 - \int d\mu(\xi^N) \xi^\top \frac{1}{N} \mathbf{s}^\top \mathbf{t} \xi \\ &= 1 - \int d\mu(\xi) \xi^\top R \xi \\ &= 1 - R \int d\mu(\xi) \xi^\top \xi, \end{aligned} \quad (95)$$

where this holds because the elements of  $\xi$  are i.i.d. Since  $d\mu(\xi^N)$  is a measure over a (multi-variate) Gaussian distribution, The data vectors  $\xi$  are normalized (See Section A.2.1) such that the second term on the R.H.S. is unity and we recover (i.e., Eqn. 30)

$$\epsilon(R) = \langle \Delta \mathbf{E}_{\ell_2}(\mathbf{s}, \mathbf{t}, \xi^n) \rangle_{\xi^N} = 1 - R. \quad (96)$$

In traditional **StatMech** (e.g., [?]), one is interested in how the *Total Model Generalization Accuracy*  $\mathcal{E}_{gen}(R)$  depends on  $R$ . With these simple error functions, Eqn. 87 reduces to a function over  $R$ , and the Average ST Generalization Error  $\mathcal{E}_{gen}^{ST}(R)$  is then obtained by averaging over the Students

$$\bar{\mathcal{E}}_{gen}^{ST}(R) = \langle \epsilon(R) \rangle_{\mathbf{s}}^\beta = \left\langle 1 - \langle \eta(\mathbf{s}, \mathbf{t}, \xi) \rangle_{\xi^N} \right\rangle_{\mathbf{s}}^\beta = \langle 1 - \eta(\mathbf{s}, \mathbf{t}) \rangle_{\mathbf{s}}^\beta = \left\langle 1 - \frac{1}{N} \mathbf{s}^\top \mathbf{t} \right\rangle_{\mathbf{s}}^\beta = \langle (1 - R) \rangle_{\mathbf{s}}^\beta, \quad (97)$$

where  $\langle \dots \rangle_{\mathbf{s}}^\beta$  is a Thermal Average over the Student weight vector  $\mathbf{s}$ .

The Model Quality for the ST Perceptron,  $\bar{\mathcal{Q}}^{ST}$  is just the Average Generalization Accuracy, so we can write

$$\bar{\mathcal{Q}}^{ST} := 1 - \bar{\mathcal{E}}_{gen}^{ST}(R) = \langle 1 - \epsilon(R) \rangle_{\mathbf{s}}^\beta = \left\langle \langle \eta(\mathbf{s}, \mathbf{t}, \xi) \rangle_{\xi^N} \right\rangle_{\mathbf{s}}^\beta = \langle \eta(\mathbf{s}, \mathbf{t}) \rangle_{\mathbf{s}}^\beta = \left\langle \frac{1}{N} \mathbf{s}^\top \mathbf{t} \right\rangle_{\mathbf{s}}^\beta = \langle R \rangle_{\mathbf{s}}^\beta. \quad (98)$$

Eqn. 98 is the starting point for deriving a **SemiEmpirical** theory for the **WeightWatcher** quality metrics (**Alpha, AlphaHat**); see Section 5.1. To generalize this expression, we will start with the Self-Overlap  $\eta(\mathbf{S}, \mathbf{T}, \xi)$  for a Multi-Layer Perceptron (MLP3) in Section 5.

Before doing this, however, we note that we can obtain this expression for  $\mathcal{E}_{gen}^{ST}$  by defining the Annealed Hamiltonian  $H_{hT}^{an}(R)$ , at high-Temperature, as in Section 4.2, Eqn. 50. Indeed, it is really  $H_{hT}^{an}(R) = \epsilon(R)$  that we must generalize to the matrix case, which we do (using a technique similar to a Replica calculation, but still in the AA). For more details, see Appendix A.2. (In particular, doing this allows us to define the normalization needed later for the **TRACE-LOG** condition).

Quantity	Traditional SMOG	Linear Perceptron in Traditional SMOG	Matrix Generalization for SETOL
Total Data Model Error	$\Delta \mathbf{E}_{\mathcal{L}}(\mathbf{w}, \boldsymbol{\xi}^n)$ (31)	$\Delta \mathbf{E}_{\mathcal{L}}(\mathbf{s}, \mathbf{t}, \boldsymbol{\xi})$ (90)	$\Delta \mathbf{E}_{\ell_2}(\mathbf{S}, \mathbf{T}, \boldsymbol{\xi})$ (101)
Annealed Hamiltonian (Data-Averaged Error)	$H_{hT}^{an} = \epsilon(\mathbf{w})$ (30) (AA, at high-T)	$H_{hT}^{an}(R) = \epsilon(\mathbf{s}, \mathbf{t}) = 1 - R$ (96) (and at large-N)	$H_{hT}^{an}(\mathbf{R}) = \mathbf{I}_M - \mathbf{R}$ (193) (only for a layer)
Self-Overlap	$\eta(\mathbf{w}) = 1 - \epsilon(\mathbf{w})$	$\eta(\mathbf{s}, \mathbf{t}) = \frac{1}{N} \mathbf{s}^\top \mathbf{t}$ (93)	$\eta(\mathbf{S}, \mathbf{T}) = \frac{1}{N} \mathbf{S}^\top \mathbf{T}$ (106)
Model Quality in terms of Layer Quality	$\bar{\mathcal{Q}} := 1 - \bar{\mathcal{E}}_{gen}$	$\bar{\mathcal{Q}}^{ST} := 1 - \bar{\mathcal{E}}_{gen}^{ST}$ (59)	$\bar{\mathcal{Q}}^{NN} := 1 - \bar{\mathcal{E}}_{gen}^{NN}$ (59) $\bar{\mathcal{Q}}^{NN} := \prod_L \bar{\mathcal{Q}}_L^{NN}$

Table 4: Summary of key quantities compared across traditional SMOG models, the Student-Teacher (ST) Linear Perceptron–in the Annealed Approximation (AA) and at high-Temperature (high-T) and at large- $N$ , and the matrix-generalized forms as the starting point to frame SETOL. The Total Data Model Error represents the difference between the model and its labels for the ST model between the Student and Teacher predictions. The Annealed Hamiltonian is the Energy function for this Error after it is averaged over the model for the training data (an  $n$  or  $N$ -dimensional i.i.d. Gaussian model, i.e.,  $\boldsymbol{\xi}^n$  or  $\boldsymbol{\xi}^N$ ). In the AA, the Annealed Hamiltonian is equal to the Annealed Error Potential. For the ST model, this is one minus the overlap,  $(1 - R)$ ; for the SETOL, this is the ( $M$ -dimensional) identity minus the overlap operator/matrix,  $\mathbf{I}_M - \mathbf{R}$ . The Self-Overlap  $\eta(\dots)$  is used to describe the Accuracy (as opposed to the Error) for both the ST model and its matrix-generalized form. Finally, the different forms of the Quality are defined. Generally speaking, the Quality  $\bar{\mathcal{Q}}$  is an approximation to some measure of 1 minus the Average Generalization Error,  $(\bar{\mathcal{Q}} := 1 - \bar{\mathcal{E}}_{gen})$  (in the AA, at high-T, at large-N, and with whatever else approximations are applied). For the ST model, having just 1 layer, the Model Quality and the Layer Quality are the same, and denoted  $\bar{\mathcal{Q}}^{ST}$ . For SETOL, the Model Quality  $\bar{\mathcal{Q}}^{NN}$  is a product of individual Layer Qualities  $\bar{\mathcal{Q}}_L^{NN}$ . (Note that the final SETOL Layer Quality  $\bar{\mathcal{Q}}$  is defined in terms of the Layer Quality-Squared  $\bar{\mathcal{Q}}^2$ , and the starting point for this is expressed with the Layer Quality-Squared Hamiltonian  $\mathbf{H}_{\bar{\mathcal{Q}}^2} = \mathbf{R}^\top \mathbf{R}$ ).

## 5 Semi-Empirical Theory of the HTSR Phenomenology

In this section, we present the main technical elements of our Semi-Empirical Theory of Deep Learning (SETOL). Our goal is to explain and, where possible, derive the HTSR PL metrics Alpha ( $\alpha$ ) and AlphaHat  $\hat{\alpha}$  from first principles, and, in doing so, also present the TRACE-LOG condition and newly proposed WeightWatcher DetX metric. To do this, we introduce a Matrix Generalization of the Student-Teacher model for a Linear Perceptron (See Section 4.3.2), adapted here for a (3-layer) Multi-Layer Perceptron (MLP3). We seek a theory for the Layer Quality  $\bar{\mathcal{Q}} = \bar{\mathcal{Q}}_L^{NN}$  of a NN, where this Layer Quality now corresponds to the (approximate) contribution each layer makes to the total generalization accuracy, or total Quality  $\bar{\mathcal{Q}}^{NN}$ . For technical reasons, we actually seek a formal expression(s) for the Layer Quality-Squared,  $\bar{\mathcal{Q}}^2 \approx \bar{\mathcal{Q}}^2$ . We say that the SETOL is Semi-Empirical because the final result  $\bar{\mathcal{Q}}^2$  is expressed directly in terms of the empirically observable spectral properties of the Teacher layer weight matrix  $\mathbf{T} = \mathbf{W}$ .

**5.1 Matrix Generalization of the ST Model.** Section 5.1 generalizes classical StatMech vector-based ST model of Section 4.3 to obtain a Layer Quality for a single layer in a NN. It starts by first formulating the learning problem for the NN generalization accuracy or quality,  $\bar{\mathcal{Q}}^{NN}$ , of a 3-layer MLP (MLP3). We then replace vectors with  $N \times M$  matrices  $\mathbf{s}, \mathbf{t} \rightarrow \mathbf{S}, \mathbf{T}$ , and obtain an expression for the NN Self-Overlap  $\eta(\mathbf{S}, \mathbf{T}, \boldsymbol{\xi})$ , which then gives a

matrix-generalized overlap operator  $\mathbf{R} := \langle \eta(\mathbf{S}, \mathbf{T}, \boldsymbol{\xi}) \rangle_{\boldsymbol{\xi}^N}^{-1} = \frac{1}{N} \mathbf{S}^\top \mathbf{T}$ . This can be related to a single-layer matrix-generalization of the ST Annealed Hamiltonian, presented in Appendix A.2,  $H_{hT}^{an} := M - \mathbf{R}$ , where, importantly, the scalar overlap  $R$  is now a matrix  $\mathbf{R}$  of  $M$  adjustable parameters.

**5.2 The Layer Quality-Squared  $\bar{Q}^2$**  Section 5.2 presents the expression for NN Layer Quality-Squared  $\bar{Q}^2$ . Following the ST analogy, we define a Thermal Average over possible Student weight matrices  $\mathbf{S}$  for the matrix overlap, giving  $\bar{Q}_L^{NN} := \langle H_{hT}^{an} \rangle_{\mathbf{S}}^\beta = \langle \mathbf{R} \rangle_{\mathbf{S}}^\beta$ . For technical reasons, however, we actually seek the (approximate) Layer Quality-Squared,  $Q^2 \approx \bar{Q}^2$ , defined as  $\bar{Q}^2 := \langle \mathbf{R}^\top \mathbf{R} \rangle_{\mathbf{S}}^\beta$ . To evaluate  $\bar{Q}^2$ , rather than sampling all random Student matrices  $\mathbf{S}$  directly, we switch measures to the Student correlation matrices  $\mathbf{A}_2 = \frac{1}{N} \mathbf{S} \mathbf{S}^\top$ . Importantly, we argue that the measures  $d\mu(\mathbf{A}_1) \leftrightarrow d\mu(\mathbf{A}_2)$ , can be interchanged for our purposes, making them effectively equivalent. This reparametrization leads us to an integral of the HCIZ type (as in Eqn. 70) which, as shown by Tanaka [?, ?].

Then, we introduce the Effective Correlation Space (ECS), and two key approximations, the TRACE-LOG Condition and the Independent Fluctuation Approximation (IFA). The TRACE-LOG condition states that the determinant of the (effective) Student correlation matrix is unity,  $\det(\tilde{\mathbf{A}}) = 1$ . Critically, this condition can be tested empirically by assuming the (effective) Teacher correlation matrix also follows the TRACE-LOG condition,  $\det(\tilde{\mathbf{X}}) = 1$ , and this is a key result of this work. Finally, we impose the IFA (described below) because it is necessary for the final result.

**5.3 The Large- $N$  limit** Section 5.3 presents the core result, (as in Eqn. 15), a closed-form or semi-analytic expressions for the Layer Quality-Squared  $\bar{Q}^2$  formed in the large- $N$  limit. Restricted to the ECS, and under the TRACE-LOG condition and the IFA, our HCIZ integral for  $\bar{Q}^2$  becomes tractable at large- $N$ , giving an expression that can be parameterized in terms of  $\tilde{M}$  eigenvalues  $\tilde{\lambda}$  of the Teacher correlation matrix restricted to the ECS  $\tilde{\mathbf{X}}$ . In doing this, the  $\tilde{M}$  Teacher eigenvalues are treated as experimental observables, and become the effective Semi-Empirical parameters (i.e.,  $\alpha$ ,  $\lambda_{max}$ ) of our SETOL.

**5.4 Modeling the Heavy-Tailed R-transform** Section 5.4 presents several models of different R-transforms. Evaluating  $\bar{Q}^2$  requires evaluating selecting an R-transform  $R(z)$  for the Teacher ESD, and also ensure that it is analytic and single-valued on the domain of interest—the ECS and/or tail of the ESD. We examine three possible models for  $R(z)$ : (i) the *Spikes-only model*, (ii) the *Inverse Wishart* (IW) model, and (iii) the *Levy-Wigner* (LW) model. First, as a trivial case, the tail of ESD can be treated a collection of spikes, and the ESD is simply as a sum of Dirac delta functions; in this case,  $\bar{Q}^2$  becomes a “Tail Norm” or “Trace Norm” ; check this When the layer is Ideal, i.e.,  $\alpha \sim 2$  and  $\det(\tilde{\mathbf{X}}) \sim 1$ , one can use *Inverse Wishart* (IW) model. As required, the IW R-transform contains a *branch cut* in the complex plane which aligns with the start of the ECS/Power Law tail. Finally, Using the Levy-Wigner model, one can (at least formally) derive the HTSR AlphaHat metric.

These core elements form a bridge between well-established empirical properties of large-scale NNs and a tractable ST-based theory. In the subsequent sections, we formalize the key steps: (i) setting up the matrix-based ST problem, (ii) defining our HCIZ integrals over restricted correlation matrices (ECS), and (iii) analyzing the resulting *Layer Quality* (or *Quality-Squared*) expressions in the large- $N$  limit.

## 5.1 Multi-Layer Setup: MLP3

In this section, we describe the matrix generalization of the ST model of the Linear Perceptron; and, in particular, a matrix-generalized version of the key quantities we derived in Section 4.3.2.

**A simple model.** Consider a simple NN with three layers (two hidden and an output), i.e., a three-layer Multi-Layer Perceptron, denoted as the *MLP3* model. (This is a *very* simple model of a modern NN with hundreds of layers and complex internal structure.)

Ignoring the bias terms, *Without Loss of Generality*, (WLOG), the NN outputs  $E_{1\mu}, E_{2\mu}, E_{3\mu}$  for each layer, as defined in Eqn. 1, are given by:

$$\begin{aligned} E_{1\mu} &:= \frac{1}{\sqrt{N_1}} h(\mathbf{W}_1^\top \boldsymbol{\xi}_\mu), \\ E_{2\mu} &:= \frac{1}{\sqrt{N_2}} h(\mathbf{W}_2^\top \mathbf{Z}_{1\mu}), \\ y_\mu &:= E_{3\mu} = \frac{1}{\sqrt{N_3}} h(\mathbf{W}_3^\top \mathbf{Z}_{2\mu}), \end{aligned} \quad (99)$$

where  $h$  is a general function or functional, denoting either a non-linear activation or a more complex layer structure, such as a CNN or an RNN. We can consider  $h(\cdot)$  to be an (unspecified) activation function.

Let us specify the ST error, or Energy difference, specifically in terms of the L2 or RMSE loss:

$$\Delta \mathbf{E}_{\ell_2}(\mathbf{S}, \mathbf{T}, \boldsymbol{\xi}) = \frac{1}{2} \sum_{\mu=1}^N (y_\mu^S - y_\mu^T)^2. \quad (100)$$

We now start to develop the self-overlap  $\eta$  using  $\mathbf{S} - \mathbf{T}$ :

### 5.1.1 Data-Dependent Multi-Layer ST Self-Overlap ( $\eta(\mathbf{S}, \mathbf{T})$ )

It is convenient to rewrite  $\Delta \hat{\mathbf{E}}$  in Eqn. 100 as:

$$\Delta \mathbf{E}_{\ell_2}(\mathbf{S}, \mathbf{T}, \boldsymbol{\xi}) := \frac{1}{2} \text{Tr}[(y_\mu^S - y_\mu^T)^\top (y_\mu^S - y_\mu^T)] = N - \text{Tr}[(\mathbf{y}^S)^\top \mathbf{y}^T] = N(1 - \eta([\mathbf{S}_l, \mathbf{T}_l], \boldsymbol{\xi})) \quad (101)$$

where the Self-Overlap  $\eta([\mathbf{S}_l, \mathbf{T}_l], \boldsymbol{\xi})$  is of the same form as the (vector) Linear Perceptron (in Eqn. 90). Note that  $\eta([\mathbf{S}_l, \mathbf{T}_l], \boldsymbol{\xi})$  depends on the (model) data  $\boldsymbol{\xi}$  because we have not evaluated the expected value  $\langle \dots \rangle_{\boldsymbol{\xi}}$  yet.

Using the general expression from Eqn. 99 for the action of the NN on the input data  $\boldsymbol{\xi}$ , we can write the formal expression of the ST error for the simple MLP3 model as:

$$\begin{aligned} \langle \eta([\mathbf{S}_l, \mathbf{T}_l], \boldsymbol{\xi}) \rangle_{\bar{\boldsymbol{\xi}}} &= \langle \eta(\mathbf{S}_1, \mathbf{S}_2, \mathbf{S}_3, \mathbf{T}_1, \mathbf{T}_2, \mathbf{T}_3, \boldsymbol{\xi}) \rangle_{\bar{\boldsymbol{\xi}}} \\ &:= \text{Tr} \left[ \frac{1}{\sqrt{N_3}} h \left( \mathbf{S}_3^\top \frac{1}{\sqrt{N_2}} h \left( \mathbf{S}_2^\top \frac{1}{\sqrt{N_1}} h(\mathbf{S}_1^\top \boldsymbol{\xi}) \right) \right) \right]^\top \times \frac{1}{\sqrt{N_3}} h \left( \mathbf{T}_3^\top \frac{1}{\sqrt{N_2}} h \left( \mathbf{T}_2^\top \frac{1}{\sqrt{N_1}} h(\mathbf{T}_1^\top \boldsymbol{\xi}) \right) \right) \right]. \end{aligned} \quad (102)$$

So far, we have not used any particular assumption on the form of the NN or the data, other than that the layer structure used to write the explicit expression for the form eventually needed,  $\eta(\mathbf{S}, \mathbf{T})$ , a single layer Self-Overlap. As a next step, we show which assumptions are needed in order to reformulate the setup as an effectively a single layer linear model for a NN.

### 5.1.2 A Single Layer Matrix Model

Following others in the literature [?], and for simplicity, one can restrict to the simplifying case that the function  $h(x)$  is the identity function. To evaluate Eqn. 102, there are three possibilities. First, we can multiply all the matrices together, and treat a multi-layer NN effectively as a single layer. Under this assumption, Eqn. 102 simplifies to

$$\langle \eta([\mathbf{S}_l, \mathbf{T}_l], \boldsymbol{\xi}) \rangle_{\bar{\boldsymbol{\xi}}^N} = \frac{1}{N_3 N_2 N_1} \text{Tr} [\boldsymbol{\xi}^\top \mathbf{S}_1 \mathbf{S}_2 \mathbf{S}_3 \mathbf{T}_3^\top \mathbf{T}_2^\top \mathbf{T}_1^\top \boldsymbol{\xi}]. \quad (103)$$

While this is possible, it would not lead to layer-by-layer insights (as HTSR-based approaches do). Second, we could attempt to expand Eqn. 103 into inter- and intra-layer terms, which we could readily do if the  $S$  and  $T$  matrices were square and the same shape, and then apply Wicks theorem:

$$\eta([\mathbf{S}_l, \mathbf{T}_l]) \approx \prod_{l=1}^L \frac{1}{N_l} \langle \eta(\mathbf{S}_l, \mathbf{T}_l, \boldsymbol{\xi}) \rangle_{\bar{\boldsymbol{\xi}}^N} = \prod_{l=1}^L \frac{1}{N_l} \text{Tr} [\boldsymbol{\xi}^\top \mathbf{S}_l^\top \mathbf{T}_l \boldsymbol{\xi}] + \text{intra-layer cross terms}. \quad (104)$$

However, these matrices are not square, and we don't know how to express the intra-layer cross terms. Finally, we can simply assume that the individual layers are statistically independent, in which case we can treat each layer independently. By ignoring the intra-layer cross-terms, let us write the single-layer self-overlap  $\eta(\mathbf{S}_l, \mathbf{T}_l, \boldsymbol{\xi})$  as:

$$\eta(\mathbf{S}_l, \mathbf{T}_l, \boldsymbol{\xi}) = \langle \eta(\mathbf{S}_l, \mathbf{T}_l, \boldsymbol{\xi}) \rangle_{\bar{\boldsymbol{\xi}}^N} \rightarrow \frac{1}{N} \text{Tr} [\boldsymbol{\xi}^\top \mathbf{S}_l^\top \mathbf{T}_l \boldsymbol{\xi}]. \quad (105)$$

This third approach is the one we will adopt. Moving forward, we will drop the layer subscript,  $l$ , and we will consider a SETOL as a single-layer theory.

### 5.1.3 The Matrix-Generalized ST Overlap ( $\eta(\mathbf{S}, \mathbf{T})$ )

We can now relate the Self-Overlap for the NN layer  $\eta(\boldsymbol{\xi})_l$  in Eqn. 105 as a matrix-generalized form of the ST Annealed Error Potential  $\epsilon(R) = 1 - R$ . Since we can interpret the Trace as an expected value over the model data  $\boldsymbol{\xi}^N$ , this gives the desired

$$\eta(\mathbf{S}, \mathbf{T}) := \frac{1}{N} \text{Tr} [\mathbf{S}^\top \mathbf{T}] \quad (106)$$

This is the matrix generalized form of Eqn. 98 for the Layer Quality.

## 5.2 Quality Metrics of an Individual Layer as an HCIZ Integral

In this subsection, we describe how to generalize the (Thermal) average over the Students  $\langle \dots \rangle_{\mathbf{S}}^\beta$  to an integral over random Student matrices,  $\langle \dots \rangle_{\mathbf{S}}^\beta$ , called an HCIZ integral.

### 5.2.1 A Generating Function Approach to Average Quality-Squared of a Layer

For our matrix generalization, we need to express the Layer Quality  $\bar{Q}$  in terms of the data-averaged Self-Overlap in Eqn. ?? for a individual layer.

- **Student-Teacher Overlap  $\mathbf{R}$**  For the vector-based Perceptron ST model, the data-averaged Self-Overlap appears in the expression for the Layer Quality in Eqn. 98, and is just the ST vector overlap  $R = \frac{1}{N} \mathbf{s}^\top \mathbf{t}$ . For our matrix generalization, we can define

$$\mathbf{R} = \frac{1}{N} \mathbf{S}^T \mathbf{T}. \quad (107)$$

For the vector-based ST model, the ST Quality  $\bar{Q}^{ST}$  in Eqn. 98 is expressed as the Thermal Average  $\bar{Q}^{ST}(R) = \langle R \rangle_{\mathbf{S}}^\beta$ . Here, we want something similar.

- **Model and Layer Qualities**  $\bar{Q}^{NN}$ ,  $\bar{Q}_L^{NN}$  For our matrix generalization, the Model Quality  $\bar{Q}^{NN}$ , as explained in Subsection 2.3, Eqn. 7, will be a product of individual NN Layer Qualities  $\bar{Q}_L^{NN}$ , and, as in Eqn. 59 approximates the total NN Average Generalization Accuracy ( $1 - \bar{\mathcal{E}}_{gen}^{NN}$ ):

$$\bar{Q}^{NN} := \prod_L \bar{Q}_L^{NN} \approx 1 - \bar{\mathcal{E}}_{gen}^{NN} \quad (108)$$

The individual  $\bar{Q}_L^{NN}$  expresses the contribution that layer makes to the approximate total NN Average Generalization Accuracy.

- **Layer Quality-Squared**  $\bar{Q}^2$  For technical convenience, however, rather than compute the NN Layer Quality  $\bar{Q}_L^{NN}$  directly, we will work with the *Average Layer Quality-Squared*, defined as

$$\bar{Q}^2 := \langle \mathbf{R}^\top \mathbf{R} \rangle_{\mathbf{S}}^\beta \quad (109)$$

where  $\langle \dots \rangle_{\mathbf{S}}^\beta$  is now a Thermal Average over Student weight matrices  $\mathbf{S}$ —an HCIZ integral. This choice means that final Layer Quality  $\bar{Q}$  we use approximates what would be the matrix-generalized NN Layer Quality  $\bar{Q}_L^{NN}$  (above) as

$$\begin{aligned} \bar{Q} &:= \sqrt{\bar{Q}^2} = \sqrt{\langle \mathbf{R}^\top \mathbf{R} \rangle_{\mathbf{S}}^\beta} \\ &\approx \langle \sqrt{\mathbf{R}^\top \mathbf{R}} \rangle_{\mathbf{S}}^\beta \\ &\approx \langle \mathbf{R} \rangle_{\mathbf{S}}^\beta \\ &= \bar{Q}_L^{NN} \end{aligned} \quad (110)$$

- **Overlap Squared** The Overlap operator (squared)  $\text{Tr}[\mathbf{R}^\top \mathbf{R}]$  is defined in terms of Eqn. 115 such that we can

$$\text{Tr}[\mathbf{R}^\top \mathbf{R}] := \frac{1}{N^2} \text{Tr}[\mathbf{T}^\top \mathbf{S} \mathbf{S}^\top \mathbf{T}] = \frac{1}{N} \text{Tr}[\mathbf{T}^\top \mathbf{A}_2 \mathbf{T}]. \quad (111)$$

This choice places  $\bar{Q}^2$  in the form of the HCIZ integral, as in Eqn. 70. See Appendix A.6 for a detailed discussion of why we choose  $\mathbf{A} = \mathbf{A}_2$  here.

- **Generating Function** For the vector-based ST model, we could compute  $\bar{Q}^{ST}$  using a generating function,  $\beta \mathbf{\Gamma}_{\bar{Q}}^{ST}$ , For our matrix generalization, we compute  $\bar{Q}$  from a *Layer Quality-Squared Generating Function*  $\beta \mathbf{\Gamma}_{\bar{Q}^2}^{IZ}$ , given as

$$\beta \mathbf{\Gamma}_{\bar{Q}^2}^{IZ} := \ln \int d\mu(\mathbf{S}) \exp[N\beta \text{Tr}[\mathbf{R}^\top \mathbf{R}]]. \quad (112)$$

See Appendix A.3 for the derivation of Eqn. 112 (and recall the discussion in Section 4.2) . .

We can not evaluate Eqn. 112 directly; but we will be able to evaluate it if we transform it into an HCIZ integral (as in Eqn. 70). To do this, however, requires a trick which will allow us to work with different forms of the Student  $\mathbf{A}$  (and Teacher  $\mathbf{X}$ ) Correlation matrices.



1771 **From weight matrices to Correlation matrices.** To evaluate  $\bar{\mathcal{Q}}^2$  in terms of derivatives of  
 1772 Eqn. 112, we need to introduce the change of measure:

$$d\mu(\mathbf{S}) \rightarrow d\mu(\mathbf{A}), \quad (113)$$

1773 and then restrict  $d\mu(\mathbf{A})$  to resemble just the generalizing eigencomponents of the Teacher correlation  
 1774 matrix  $\mathbf{X}$ . We have two choices for the Student Correlation Matrix  $\mathbf{A}$ , call them  $\mathbf{A}_1$  and  $\mathbf{A}_2$ ,  
 1775 defined as:

$$\mathbf{A}_1 := \frac{1}{N} \mathbf{S}^\top \mathbf{S} \quad (\text{which is } M \times M) \quad (114)$$

$$\mathbf{A}_2 := \frac{1}{N} \mathbf{S} \mathbf{S}^\top \quad (\text{which is } N \times N). \quad (115)$$

1776 Note that  $\frac{1}{N}$  is the correct scaling on each of these. Eqn. 114 is consistent with our definition of the  
 1777 layer Correlation Matrix, and we use it as the starting point below to derive the Volume-Preserving  
 1778 TRACE-LOG condition (Appendix A.4). Eqn. 115 is consistent with Tanaka, which requires  $\mathbf{A}$  be  
 1779  $N \times N$ , but the we still need the a *Duality of Measures* to rederive this (Appendix A.6).[?, ?]

1780 **Duality of measures.** For either form of  $\mathbf{A}$ , the measure  $d\mu(\mathbf{A})$  is the same because we will  
 1781 restrict the measures to the ECS space of non-zero eigenvalues ( $\lambda_i \gg 0$ ). We note that  $\mathbf{A}_1$  and  $\mathbf{A}_2$   
 1782 have the same eigenvalues  $\lambda_i$ , or ESD, upto the additional zero eigenvalues ( $\lambda_i = 0$ ) in the null  
 1783 space of  $\mathbf{A}_2$ . Consequently, both forms of  $\mathbf{A}$  have the same  $N - M$  non-zero eigenvalues, the same  
 1784 non-zero part of the ESD ( $\rho(\lambda) \gg 0$ ), and the same Trace ( $\text{Tr}[\mathbf{A}_1] = \text{Tr}[\mathbf{A}_2]$ ). In the large- $N$   
 1785 approximation, the ESD of (either form of)  $\mathbf{A}$ ,  $\rho_{\mathbf{A}}^\infty(\lambda)$ , becomes continuous (and bounded), but  
 1786 remains zero in the null space. Consequently, when integrating over the eigenvalues  $\lambda$ , one can  
 1787 interchange  $\mathbf{A}_1$  with  $\mathbf{A}_2$ , such that

$$\int d\mu(\mathbf{A}) [\dots] \leftrightarrow \int d\lambda [\dots] \rho_{\mathbf{A}_1}^\infty(\lambda) \leftrightarrow \int d\lambda [\dots] \rho_{\mathbf{A}_2}^\infty(\lambda) \quad (116)$$

1788 This equivalence will be essential both to derive the TRACE-LOG condition (Section A.4), and to  
 1789 (re)derive the core result by Tanaka (Section A.6). Additionally and WLOG, we may occasionally  
 1790 denote the Student Correlation matrix as  $\hat{\mathbf{A}}$  instead of explicitly using  $\mathbf{A}_1$  and  $\mathbf{A}_2$ .

### 1791 5.2.2 Evaluating the Average Quality (Squared) Generating Function

1792 We can write the generating function  $\beta \Gamma_{\bar{\mathcal{Q}}^2}^{IZ}$  in Eqn. 112 in terms of  $\mathbf{A}_2$ , giving, as in Eqn. 12,

$$\beta \Gamma_{\bar{\mathcal{Q}}^2}^{IZ} = \ln \int d\mu(\mathbf{S}) e^{N\beta \text{Tr}[\frac{1}{N} \mathbf{T}^\top \mathbf{A}_2 \mathbf{T}]} \quad (117)$$

1793 To recast Eqn. 117 as an HCIZ integral, as in Eqn. 70, we must perform change of measure, from  
 1794  $N \times M$  Student weight matrices  $\mathbf{S}$  to  $N \times N$  Student Correlation matrices  $\mathbf{A}$ , as in Eqn. 113.

1795 When we perform the change of measure on Eqn. 117, we obtain the following expression:

$$\begin{aligned} \beta \Gamma_{\bar{\mathcal{Q}}^2}^{IZ} &\approx \ln \int d\mu(\mathbf{A}) e^{N\beta \text{Tr}[\frac{1}{N} \mathbf{T}^\top \mathbf{A}_2 \mathbf{T}]} e^{\frac{N}{2} \ln(\det(\mathbf{A}_1))} \\ &= \ln \left\langle e^{N\beta \text{Tr}[\frac{1}{N} \mathbf{T}^\top \mathbf{A}_2 \mathbf{T}]} e^{\frac{N}{2} \ln(\det(\mathbf{A}_1))} \right\rangle_{\mathbf{A}} \end{aligned} \quad (118)$$

1796 where the latter expresses the former in Bra-Ket notation. This expression is derived in  
 1797 Appendix A.4 (see Eqn. 236): it contains the original overlap term that depends  $\mathbf{A}_2$  as well as a  
 1798 new term from the transformation that depends on  $\det(\mathbf{A}_1)$  that is not yet defined.

### 5.2.3 The Effective Correlation Space (ECS)

Currently, Eqn. 118 is a “formal” expression, and we have not identified and/or justified its realm of applicability. Fortunately, we have empirical evidence to suggest that this can be made “physically” meaningful, by “restricting” the integral to the tail of the ESD.

Prior work on HTSR theory indicates that the generalizing parts of a layer weight matrix concentrate into  $\rho_{tail}^{emp}(\lambda)$ , the tail of the ESD, as the ESD becomes more Power Law (PL), and layer PL exponent  $\alpha \rightarrow 2$  (from above) [?, ?, ?, ?, ?]. This suggests the integral in Eqn. 9 should instead average over a low rank subspace spanned *only by* the generalizing eigen-components of the student layer correlation matrix,  $\mathbf{A}_1 = \frac{1}{N} \mathbf{S}^\top \mathbf{S}$  (or, equivalently,  $\mathbf{A}_2$ ). We call this subspace the Effective Correlation Space (ECS).

Let  $\tilde{\mathbf{A}}$  be the matrix spanned by the largest  $\tilde{M}$  eigen-components  $\mathbf{A}$  (in either form,  $\mathbf{A}_1$  or  $\mathbf{A}_2$ ).

$$\tilde{\mathbf{A}} := \mathbf{P}_{ecs} \mathbf{A}, \quad \mathbf{P}_{ecs} := \sum_{i=1}^{\tilde{M}} |\lambda_i\rangle \langle \lambda_i| \quad (119)$$

where  $\mathbf{P}_{ecs}$  is the projection operator onto the subspace spanned by the eigenvector  $|\lambda_i\rangle$  associated with the eigenvalue  $\lambda_i$  of  $\mathbf{A}_1$  or, equivalently,  $\mathbf{A}_2$ . We denote the corresponding Student correlation matrices with tilde, as follows:

$$\begin{aligned} \mathbf{A}_1 &\rightarrow \tilde{\mathbf{A}}_1 \text{ such that } d\mu(\mathbf{A}_1) \rightarrow d\mu(\tilde{\mathbf{A}}_1) \\ \mathbf{A}_2 &\rightarrow \tilde{\mathbf{A}}_2 \text{ such that } d\mu(\mathbf{A}_2) \rightarrow d\mu(\tilde{\mathbf{A}}_2) \end{aligned} \quad (120)$$

where now the matrices  $\tilde{\mathbf{A}}_1$  and  $\tilde{\mathbf{A}}_2$  are restricted to the ECS, and the measure  $d\mu(\tilde{\mathbf{A}}_1)$  is similarly restricted. See Appendix A.4 for a more detailed discussion of this.

When an ESD is *Fat-Tailed* the ECS is at least as large if not larger than the PL tail of the ESD. For an ESD with  $\alpha \geq 2$ , The generalizing eigen-components, i.e.  $|\lambda_i\rangle$ , will mostly concentrate into the tail of the ESD, but some will remain in the bulk, so the ECS will be larger than the tail. In this case that  $\tilde{M} \geq M^{tail}$ , and  $\rho_{ECS}(\lambda) \supseteq \rho_{tail}(\lambda)$ , as depicted in Figure XXX. When  $\alpha = 2$ , the layer is Ideal (as in Section 3.1), in that all of the generalizing eigen-components to have now concentrated completely into the tail, so that  $\tilde{M} = M^{tail}$ , and  $\rho_{ECS}(\lambda) = \rho_{tail}(\lambda)$ . (When  $\alpha < 2$ , this suggests that the layer is overfit, and the layer may have a Correlation Trap and/or frequently also has many near-zero eigenvalues.)

This leads to the following Model Selection Rule (MSR) for the ECS:

When transforming the measure  $d\mu(\mathbf{S}) \rightarrow d\mu(\tilde{\mathbf{A}})$ , we invoke an eigenvalue cutoff rule that prescribes how to replace  $\mathbf{A}$  with a low-rank effective matrix  $\mathbf{A} \rightarrow \tilde{\mathbf{A}}$ , where the cutoff  $\tilde{\lambda} \geq \lambda_{min}^{ECS}$  is chosen so that the ECS at least contains the PL tail and, importantly, such that  $\det(\tilde{\mathbf{A}}) = \det(\mathbf{A}_1) = \det(\mathbf{A}_2)$  is well defined.

Formally, this means that when we evaluate the Quality (squared)  $\bar{Q}^2$ , Generating Function  $\beta \mathbf{\Gamma}_{\bar{Q}^2}^{IZ}$  or other relevant averages, we restrict the measure (i.e., integral or sum) to the eigencomponents in the tail of the ESD of  $\mathbf{A}$  (or  $\mathbf{X}$ , when appropriate) starting with  $\lambda_{min}^{ECS}$ . To our knowledge, this proposed MSR is completely novel.<sup>31</sup>

Restricted to the ECS, we now replace Eqn. 118 with:

$$\begin{aligned} \beta \mathbf{\Gamma}_{\bar{Q}^2}^{IZ} &= \ln \int d\mu(\tilde{\mathbf{A}}) e^{N\beta \text{Tr}[\frac{1}{N} \mathbf{T}^\top \tilde{\mathbf{A}}_2 \mathbf{T}]} e^{\frac{N}{2} \ln(\det(\tilde{\mathbf{A}}_1))} \\ &= \ln \left\langle e^{N\beta \text{Tr}[\frac{1}{N} \mathbf{T}^\top \tilde{\mathbf{A}}_2 \mathbf{T}]} e^{\frac{N}{2} \ln(\det(\tilde{\mathbf{A}}_1))} \right\rangle_{\tilde{\mathbf{A}}} \end{aligned} \quad (121)$$

<sup>31</sup>It is, however, similar in spirit to cut-offs used in some Quantum Field Theories to make the theories physically meaningful.

1834 where we have used the formal Duality of Measures,  $d\mu(\tilde{\mathbf{A}}) = d\mu(\tilde{\mathbf{A}}_1) = d\mu(\tilde{\mathbf{A}}_2)$ .<sup>32</sup>

#### 1835 5.2.4 Two Simplifying Assumptions: the IFA and TRACE-LOG Condition

1836 It now remains how to define the cutoff for the ECS space. To this, we make the following  
1837 assumptions.

- 1838 • **The Independent Fluctuation Assumption (IFA).** This condition states that the two  
1839 terms appearing in the exponential of Eqn. 121 are statistically independent:

$$\beta \Gamma_{\tilde{\mathbf{Q}}^2}^{IZ} \approx \ln \left\langle e^{N\beta \text{Tr}[\frac{1}{N} \mathbf{T}^\top \tilde{\mathbf{A}}_2 \mathbf{T}]} \right\rangle_{\tilde{\mathbf{A}}} \left\langle e^{\frac{N}{2} \ln(\det(\tilde{\mathbf{A}}_1))} \right\rangle_{\tilde{\mathbf{A}}}. \quad (122)$$

- 1840 • **The Trace Log Condition (TRACE-LOG).** This condition states that the determinant of  
1841 the Student (and Teacher) Correlation matrix is unity, such that:

$$\det(\tilde{\mathbf{A}}) = 1 \quad \text{or} \quad \text{Tr}[\ln \tilde{\mathbf{A}}] = 0, \quad (123)$$

1842 (and with  $\mathbf{A}$  additionally normalized to  $M$ , as explained in the Appendix, Section A.2.1)  
1843 so that when we replace the measure over Student layer weight matrices  $d\mu(\mathbf{S})$  with a  
1844 measure over all Student Correlation matrices  $d\mu(\mathbf{A})$ , *restricted to the ECS*, the second term  
1845 in Eqn. 122 becomes unity  $\langle \exp[\frac{N}{2} \text{Tr}[\ln[\det \mathbf{A}_1]]] \rangle_{\mathbf{A}} = 1$ , and thus vanishes in the final  
1846 expression for  $\beta \Gamma_{\tilde{\mathbf{Q}}^2}^{IZ}$  (see Eqn. ??).

1847 At this point, the IFA is made purely for mathematical convenience, i.e., we have not demonstrated  
1848 it empirically, but it is not implausible as a statistical modeling assumption. On the other hand,  
1849 we can test the TRACE-LOG condition empirically; this is a critical part of our SETOL approach.

1850 Notably, when taking the large- $N$  approximation of  $\beta \Gamma_{\tilde{\mathbf{Q}}^2}^{IZ}$ , we effectively do this independently  
1851 in two steps. First, we apply a Saddle Point Approximation (SPA) to the second term in Eqn. 122,  
1852 leading to the TRACE-LOG condition (see Appendix A.4). Most importantly, we can test the  
1853 TRACE-LOG condition empirically, and this is another critical part and justification of our SETOL  
1854 approach. Second, when applying the result by Tanaka (Eqn. 72), we are applying an SPA to the  
1855 result, but assuming we are restricted to the ECS.

1856 **Empirical Tests of the TRACE-LOG Condition and the ECS** Since we require that the students  
1857 resemble the fixed Teacher, a reasonable estimate for the average of the Student Correlation matrix  
1858  $\tilde{\mathbf{A}}$  (either  $\tilde{\mathbf{A}}_1$  or  $\tilde{\mathbf{A}}_2$ ) (in the ECS) is the point estimate provided by the actual (known) fixed  
1859 Teacher correlation matrix  $\tilde{\mathbf{X}}$ , so that

$$\langle \det \tilde{\mathbf{A}}_1 \rangle_{\tilde{\mathbf{A}}} = \langle \det \tilde{\mathbf{A}}_2 \rangle_{\tilde{\mathbf{A}}} \simeq \det \tilde{\mathbf{X}} = \prod_t \lambda_t = 1 \quad \forall \lambda_t \in \rho_{tail}^{emp}(\lambda), \quad (124)$$

1860 for all eigenvalues  $\lambda_t$  in  $\rho_{tail}^{emp}(\lambda)$  the tail of the ESD of  $\tilde{\mathbf{X}}$ , i.e.,  $\lambda_t \geq \lambda_{min}^{ECS}$ . How should one choose  
1861  $\lambda_{min}^{ECS}$  in this expression? We already know that most NNs layers have Fat-Tailed ESDs, but if the  
1862 Teacher ESD is simply *Random-Like* or even Bulk+Spikes, then the expected value  $\langle \det(\mathbf{A}_1) \rangle$   
1863 itself of the determinant of a random full rank Student Correlation matrix  $\mathbf{A}_1$  might be well defined  
1864 and easy to estimate, and we might not even need to define the lower rank ECS. Indeed, in these  
1865 cases, the correlated eigencomponents, if they exist, may be buried in the bulk region of the ESD  
1866 and not readily identifiable. But because  $\rho_{tail}^{emp}(\lambda)$  is Fat-Tailed and Power Law (PL), this poses  
1867 some difficulty, which is why need to first define the Effective Correlation Space (ECS).

<sup>32</sup>Also, while we could replace  $\mathbf{T} \rightarrow \tilde{\mathbf{T}}$ , but to simplify the notation we do not do this.

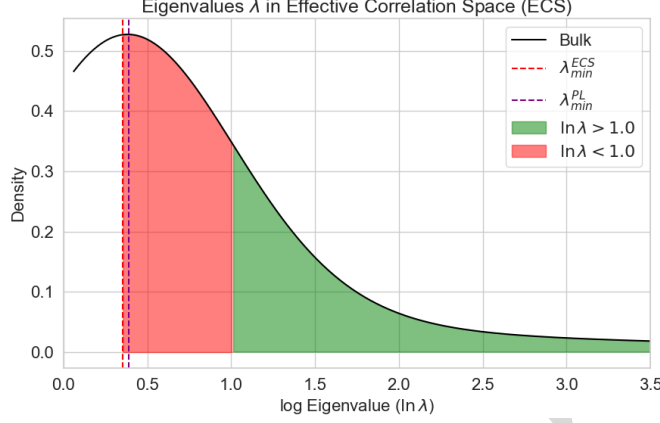


Figure 10: The image depicts a typical Empirical Spectral Density (ESD) of a layer correlation matrix  $\mathbf{X}$ , with **WeightWatcher** Power-Law (PL) exponent  $\alpha = 2.0$ . The green and red shaded regions depict eigenvalues  $\lambda$  in the Effective Correlation Space (ECS) of  $\tilde{\mathbf{X}}$ , defined by  $\lambda > \lambda_{min}^{ECS}$ . The x-axis displays the eigenvalues on the log scale,  $\ln \lambda$ . The vertical red line is at the start of the PL tail ( $\lambda_{min}^{PL}$ ). The purple, vertical line is at the start of the ECS tail ( $\lambda_{min}^{ECS}$ ). The green shaded region depicts those eigenvalues where  $\ln \lambda > 1.0$ , whereas the red shaded region depicts those eigenvalues where  $\ln \lambda < 1.0$ . The ECS is defined such that the volume-preserving **TRACE-LOG** condition is best satisfied, i.e.  $\sum \ln \lambda = 0$  for  $\lambda \geq \lambda_{min}^{ECS}$ .

1868 Because the Teacher ESD is most likely MHT and PL, if we choose  $\lambda_{min}^{ECS}$  too small, and the  
 1869 tail extends too far into the bulk region of the ESD, then for all practical purpose  $\det(\tilde{\mathbf{X}}) \ll 1$ . On  
 1870 the other hand, if we choose  $\lambda_{min}^{ECS}$  too large, then we only capture very large eigenvalues, and for  
 1871 all practical purposes  $\det(\tilde{\mathbf{X}}) \gg 1$ . Therefore, if we set the scale of  $\tilde{\mathbf{X}}$  appropriately, we can choose  
 1872 a  $\lambda_{min}^{ECS}$  such that  $\det(\tilde{\mathbf{X}}) = 1$ . In this case, by choosing  $\lambda_{min}^{ECS}$  appropriately, we can estimate the  
 1873 expected value of  $\langle \det(\mathbf{A}_1) \rangle$  with a empirical point estimate over the Teacher Correlation matrix,  
 1874 which is unity.

$$|\det \tilde{\mathbf{X}}| \simeq 1; \quad \text{Tr}[\ln \tilde{\mathbf{X}}] = \ln |\det \tilde{\mathbf{X}}| \simeq 0. \quad (125)$$

1875 This expression now be used in a practical calculation to define a low-rank subspace that both  
 1876 allows us to evaluate the HCIZ integral, and to identify, in principle, generalizing components of  
 1877 the layer. We also refer to Eqn. 125 the **TRACE-LOG** condition, which is technically its empirical  
 1878 form. For example, Figure 5.2.4 depicts the eigenvalues in the ECS for a typical ESD with PL  
 1879  $\alpha = 2.0$ , (normalized such that  $\text{Tr}[\mathbf{X}] = \|\mathbf{W}\|_F = N$ ). Notice that the start of the PL tail, ( $\lambda_{min}^{PL}$ ),  
 1880 is very close to the start of the ECS tail. ( $\lambda_{min}^{ECS}$ , i.e.  $\Delta \lambda_{min} := \lambda_{min}^{ECS} - \lambda_{min}^{PL} \approx 0$ ). Also, notice that  
 1881 while there are many large eigenvalues,  $\ln \lambda > 1.0$ , there are numerous small eigenvalues as well,  
 1882  $\ln \lambda < 1.0$ , such that the  $\sum \ln \lambda \approx 0$  for  $\lambda \geq \lambda_{min}^{ECS}$ . Additional plots like Figure 5.2.4, generated  
 1883 with **WeightWatcher**, can be found in Section 6, in Figure 21 as well as plots of  $\Delta \lambda_{min}$  vs. the  
 1884 **WeightWatcher** PL  $\alpha$  for several real-world examples, in Figures 22 and 23.

### 1885 5.3 Evaluating the Layer Quality ( $\bar{Q}$ ) in the Large- $N$ Limit

1886 To generate the Average Quality,  $\bar{Q}^2$ , we first take the large- $N$  limit of  $\beta \mathbf{\Gamma}_{\bar{Q}^2}^{IZ}$

$$\beta \mathbf{\Gamma}_{\bar{Q}^2, N \gg 1}^{IZ} := \lim_{N \gg 1} \beta \mathbf{\Gamma}_{\bar{Q}^2}^{IZ} = \lim_{N \gg 1} \ln \left\langle \exp N \beta \text{Tr} \left[ \frac{1}{N} \mathbf{T}^\top \tilde{\mathbf{A}}_2 \mathbf{T} \right] \right\rangle_{\tilde{\mathbf{A}}} \quad (126)$$

and then take the appropriate partial derivative, analogously to as we did for  $\bar{\mathcal{E}}_{gen}^{ST}$ ; see Section A.3 for more details. This gives (as in Eqn. 64)

$$\bar{\mathcal{Q}}^2 := \frac{1}{\beta} \frac{\partial}{\partial N} \beta \mathbf{\Gamma}_{\bar{\mathcal{Q}}^2, N \gg 1}^{IZ} \quad (127)$$

$$\approx_{\text{high-}T} \frac{1}{N} \frac{\partial}{\partial \beta} \beta \mathbf{\Gamma}_{\bar{\mathcal{Q}}^2, N \gg 1}^{IZ} \quad (128)$$

Notice that since we are at high-Temperature, it doesn't matter which partial derivative we take, and we expect both results to yield the same expression.

This HCIZ integral in Eqn. 126 can be evaluated (i.e in the large- $N$  limit) using a result by Tanaka —provided the matrix  $\tilde{\mathbf{A}}_2$  is low rank, which holds when the TRACE-LOG condition is satisfied. Thus, moving forward, we will assume an Effective Correlation Space (ECS) of rank  $\tilde{M}$ , where  $\lambda_{min}^{ECS}$  is the  $M^{th}$ -largest eigenvalue of  $\tilde{\mathbf{X}}$ , and defines the start of the ECS (and whatever branchcut is necessary to integrate  $R(z)$ ).

Tanaka's result for the ECS can be expressed as:

$$\lim_{N \gg 1} \frac{1}{N} \ln \langle \exp(\beta \text{Tr}[\mathbf{T}^\top \tilde{\mathbf{A}}_2 \mathbf{T}]) \rangle_{\tilde{\mathbf{A}}} = \beta \sum_{i=1}^{\tilde{M}} \mathcal{G}(\tilde{\lambda}_i), \quad (129)$$

where the sum now only includes the eigenvalues of  $\tilde{\mathbf{X}}$  (in the ECS),  $\beta = \frac{1}{T}$  is the Inverse-Temperature, and  $\tilde{\lambda}$  is an eigenvalue of  $\tilde{\mathbf{X}}$ , the Teacher Correlation matrix projected into the ECS space.  $\tilde{\mathbf{A}}_2$  is the  $N \times N$  form of the Student Correlation matrix, with  $N - M$  non-zero eigenvalues, and  $\mathbf{T}$  is the  $N \times M$  Teacher weight matrix (also effectively projected into the ECS, i.e.  $\mathbf{T} = \tilde{\mathbf{T}}$  here).  $\mathcal{G}(\lambda_i)$  is the Norm Generating Function, defined below.<sup>33</sup> This gives

$$\beta \mathbf{\Gamma}_{\bar{\mathcal{Q}}^2, N \gg 1}^{IZ} = N \beta \sum_{i=1}^{\tilde{M}} \mathcal{G}(\tilde{\lambda}_i), \quad (130)$$

This gives a final expression for the Average Layer Quality (Squared)  $\bar{\mathcal{Q}}^2$  as

$$\bar{\mathcal{Q}}^2 = \sum_{i=1}^{\tilde{M}} \mathcal{G}(\tilde{\lambda}_i), \quad (131)$$

Note that  $\bar{\mathcal{Q}}^2$  is independent of  $N$  and  $\beta$ , and, indeed, Eqn. 127 is an equality.

The average Quality (squared) can be expressed as a sum over Generating Functions  $\mathcal{G}(\lambda)$ , which depend only the statistical properties of the actual Teacher Correlation matrix  $\tilde{\mathbf{X}}$  (projected into the ECS). Each term in the sum,  $\mathcal{G}(\tilde{\lambda}_i)$ , takes the form

$$\mathcal{G}(\lambda) := \int_0^\lambda R_{\mathbf{A}}(z) dz \xrightarrow{\text{ECS}} \int_{\lambda_{min}^{ECS}}^\lambda R_{\tilde{\mathbf{A}}}(z) dz \quad (132)$$

where  $R_{\tilde{\mathbf{A}}}(\tilde{\lambda})$  is the R-transform from RMT, and  $\lambda_{min}^{ECS}$  is the lower bound of the ECS spectrum. Importantly, the R-transform for a Heavy-Tailed ESD may have a branchcut at or near the start of the ECS (as explained in Section 5.4), so restricting the integrand to start at  $\lambda_{min}^{ECS}$  is critical.

<sup>33</sup>We use the notation  $\langle \dots \rangle_{\tilde{\mathbf{A}}}$  for expected value and placed  $\frac{1}{N}$  on the L.H.S. to help the reader compare this to the original expressions in [?, ?]. Also, in [?, ?],  $\beta = 1/2$ , but, in fact, if one inserts  $-\beta$  as an inverse-Temperature into the final expression, it simply factors out.

Since we expect the best Students matrices to resemble the actual Teacher matrices, we expect the Student correlation matrix  $\tilde{\mathbf{A}}$  to have similar spectral properties to our actual empirical correlation matrices  $\tilde{\mathbf{X}}$ . That is, from the perspective of HTSR theory and the classification into 5+1 Phases of Training [?], we expect all the  $\tilde{\mathbf{A}}$  to be in the same phase as  $\tilde{\mathbf{X}}$  (and, in addition, to have the same PL exponent value). That is,

*We expect the R-transform of  $\tilde{\mathbf{A}}$  to have the same functional form as the R-transform of  $\tilde{\mathbf{X}}$ .*

If our (Teacher) NN weight matrix exhibits a HT PL, then the tail the ESD ( $\rho_{tail}(\lambda)$ ) of the Student and Teacher will both take the limiting form of a PL, with the same empirical variance  $\sigma^2$  and (critically) the same PL exponent  $\alpha$ :

$$\rho_{tail}[\tilde{\mathbf{A}}](\lambda) \sim \rho_{tail}[\tilde{\mathbf{X}}](\lambda) \sim \lambda^{-\alpha}. \quad (133)$$

Up until this point, our derivation of  $\bar{\mathcal{Q}}^2$  only depends on the TRACE-LOG condition, irrespective of the exact functional form of  $R(z)$ , therefore the SETOL approach tested by examining how well the TRACE-LOG condition holds for the layers in very well performing models. We do this in Section 6.3.

## 5.4 Modeling the R-Transform

In this section, we explain how to select the R-transform  $R(z)$  and evaluate the Norm Generating Function  $\mathcal{G}(\lambda)$  under different modeling assumptions. To apply SETOL, the model satisfy the TRACE-LOG condition—which occurs during the case of Ideal Learning. For most cases of NN models, the ESD are HT; and this, in practice, one usually would select  $R(x)$  that reflects this. We analyze several cases, noting their applicability to real-world scenarios. Most importantly, we derive expressions that resemble the WeightWatcher AlphaHat metric, at least formally valid for the case  $\alpha \approx 2$ .

### 5.4.1 Elementary Random Matrix Theory

We begin with some useful notions definitions from Random Matrix Theory. Using the ESD  $\rho(\lambda)$ , defined as

$$\rho(\lambda) := \frac{1}{N} \sum_i \delta(\lambda - \lambda_i), \quad (134)$$

we can express the *Greens Function* (or *Cauchy-Stieltjes transform*) by<sup>34</sup>

$$G(z) := \int d\lambda \frac{\rho(\lambda)}{z - \lambda}. \quad (135)$$

From  $G(z)$ , we can recover the ESD,  $\rho(\lambda)$ , using the inversion relation

$$\rho(\lambda) = \lim_{\epsilon \rightarrow 0+} \frac{1}{\pi} \text{IM}(C(\lambda + i\epsilon)), \quad (136)$$

where IM is the imaginary part of  $G(z)$ , and where the  $\lim_{\epsilon \rightarrow 0+}$  means to take the limit approaching from the upper half of the complex plane. The R-transform,  $R(z)$ , can be defined using the Blue function  $B(z)$

$$R(z) := B(z) - \frac{1}{z}, \quad (137)$$

---

<sup>34</sup>Please notice our naming and sign convention in Eqn. ???. Some authors equate the Greens Function  $G(z)$  with the Cauchy-Stieltjes transform, whereas we define  $C(z) = -G(z)$ .



1941 where the Blue function  $B(z)$  [?] is the functional inverse of the Greens Function  $G(z)$ ,<sup>35</sup> satisfying

$$B[G(z)] = z. \quad (138)$$

1942 By specifying the  $R(z)$  transform, we specify the complete ESD,  $\rho(\lambda)$ . Here, we are actually  
 1943 only interested in the tail of  $\rho(\lambda)$ . That is, we can given  $R(z) = R(z)_{tail} + R(z)_{bulk}$ , we only need  
 1944  $R(z) \approx R(z)_{tail}$ .

#### 1945 5.4.2 Known R-transforms and Analytic (Formal) Models

1946 There only a few known analytic results for the explicit R-transform  $R(z)$ . The ones we need are  
 1947 in Table 5. Below, we review some of them, explaining what ESD they correspond to, and what  
 1948 the resulting Norm Generating Function  $G(\lambda)$  would be if applied as a model  $R(x)$  in the SETOL  
 1949 approach.

[Note: removed *Multiplicative-Wishart* Might want to comment on it ]

Model	HTSR Universality class	R-transform $R(z)$
Discrete	Spikes	$\frac{1}{M} \sum_{i=1}^M \lambda_i$
Wishart Models		
Inverse-Wishart	HT/VHT	$\frac{\kappa - \sqrt{\kappa(\kappa - 2z)}}{z}$
Levy Wigner		
General ( $\alpha_l \neq 1$ )	VHT/HT	$a + bz^{\alpha_l - 1}$
Cauchy $\alpha_l = 2, \beta = 0$	$\alpha = 2$	$a - i\gamma$

Table 5: Known R-transforms for different matrix models. For the *Inverse-Wishart*, as given by Bun [?],  $\kappa = \frac{1}{2}(Q - 1)$  where,  $q = \frac{1}{Q} = \frac{M}{N} \leq 1$ . The *Levy-Wigner* model describes Wigner-like Square Random Matrices (as opposed to Wishart-like or Correlation Matrices), where the elements are drawn from a Levy-Stable distribution. The Levy-Stable  $R(z)$  is parameterized by a (real) shift parameter  $a$ , a complex phase factor  $b$  (that depends on 3 real parameters  $\alpha_l, \beta$ , and  $\gamma$ ), and, importantly, a PL-like tail exponent  $\alpha_l \in (0, 2)$ ; For more details, the text, see [?, ?, ?]. **For our modeling purposes here, we make the association  $\alpha_l \sim \alpha/2 - 2$ . maybe not?** (Also, for simplicity, we assume the variance  $\sigma = 1$  for all models above, where appropriate.)

1950

#### 1951 5.4.3 Discrete Model: Spikes

1952 Here, we consider modeling the HT tail ESD,  $\rho_{tail}(\lambda)$ , as a collection of discrete spikes  $\lambda_{spike}$ ,  
 1953 where  $\lambda_{spike} \geq \lambda_{min}^{ECS}$ . Here, R-transform for the ECS is sum of Dirac delta functions (as opposed to  
 1954 using the Inverse-Wishart (IW) or other model  $R(z)$ .) This lets us compute the Layer Quality  $\bar{Q}$   
 1955 in closed form in term of the Teacher weight matrix  $\mathbf{T} = \mathbf{W}$ .

<sup>35</sup>The Blue function was first introduced by Zee [?] to model, among other things, spectral broadening in quantum systems. Briefly, given a deterministic Hamiltonian matrix  $\mathbf{H}_0$ , with eigenvalues  $\lambda_i^0$ , one can model the spectral broadening of  $\lambda_i^0$  by adding a random matrix  $\mathbf{H}_1$  to  $\mathbf{H}_0$ :  $\mathbf{H} = \mathbf{H}_0 + \mathbf{H}_1$ . The resulting eigenvalues of  $\mathbf{H}$  now contain some level of randomness,  $\sigma$ , i.e.,  $\lambda = \lambda^0 + \sigma$ . To model the ESD of  $\mathbf{H}$ , one then specifies the individual R-transforms for  $\mathbf{H}_0$  and  $\mathbf{H}_1$ ; the full ESD of  $\mathbf{H}$  can then be reconstructed by adding the two R-transforms together  $R(z) = R_0(z) + R_1(z)$ . Zee also notes that  $R(z)$  is the same as the self-energy  $\Sigma(z)$  from quantum many body theory [?].

1956 Let the tail of the ESD have  $\tilde{M} = M^{tail}$  eigenvalues that define the ECS, i.e.,

$$\rho_{tail}(\lambda) = \sum_{i=1}^{\tilde{M}} \delta(\lambda - \lambda_i). \quad (139)$$

1957 The Greens Function  $G(z)$  is then

$$G(z) = \int d\lambda \frac{\rho_{tail}(\lambda)}{z - \lambda} = \sum_i \int d\lambda \frac{\delta(\lambda - \lambda_i)}{z - \lambda} = \sum_i \frac{1}{z - \lambda_i}, \quad (140)$$

1958 and the Blue function for each individual term  $i$  is  $\frac{1}{z - \lambda_i}$ , i.e.,  $B(z) = \lambda_i + \frac{1}{z}$ . Now, using the additive  
1959 property of the R-transform, we can express the total  $R(z)$  as the sum of the R-transforms for the  
1960 individual terms  $i$ , giving

$$R(z) = \sum_i \left( \lambda_i + \frac{1}{z} \right) - \frac{1}{z} = \sum_i \lambda_i. \quad (141)$$

1961 This gives the Norm Generating Function  $\mathcal{G}(\lambda)$  as

$$\begin{aligned} \mathcal{G}(\lambda) &= \int_{\lambda_{min}^{ECS}}^{\lambda} \sum_i \lambda_i d\lambda \\ &= \sum_i \lambda_i \int_{\lambda_{min}^{ECS}}^{\lambda} 1 d\lambda \\ &= \left( \sum_i \lambda_i \right) (\lambda - \lambda_{min}^{ECS}) \end{aligned} \quad (142)$$

1962 Seeing that  $\lambda_{min}^{ECS}$  is usually quite small, we make the approximation

$$\mathcal{G}(\lambda) \approx \left( \sum_i \lambda_i \right) (\lambda), \quad (143)$$

1963 which gives the Quality-Squared approximately as

$$\bar{Q}^2 = \sum_i \mathcal{G}(\lambda_i) \approx \left( \sum_i \lambda_i \right)^2. \quad (144)$$

1964 We now see that we can define  $\bar{Q} := \sqrt{\bar{Q}^2} = \sum_i \lambda_i$  is what we call a Tail Norm, the  $L1$  norm of the  
1965 tail eigenvalues.

1966 [CHECK FOR ERRORS PLEASE; Also ? not really a Trace norm since this is defined over singular  
1967 values ?]

#### 1968 5.4.4 Inverse-Wishart Model of Ideal Learning

1969 Here, we consider the Inverse-Wishart (IW) model. The IW model treats the ESD of  $\mathbf{X}^{-1}$  when  
1970 the ESD of  $\mathbf{X}$  itself is MP. As a parametric model, it can be quite effective at treating VHT and  
1971 HT (or Fat-Tailed) ESDs when the far tail decays very rapidly, like a TPL, and/or for  $\alpha \leq 4$ . To  
1972 do this, one simply considers the parameters  $\kappa$  as an adjustable parameter. It is an excellent  
1973 model for the ESD when  $\alpha = 2.0$  (and  $Q = 2$ ). Using this model, we can derive an expression for  
1974 the HTSR AlphaHat Layer Quality metric,  $\hat{\alpha} := \alpha \log_{10} \lambda_{max}$  as a leading order term in the final  
1975 expression for  $\log_{10} \bar{Q}^2$ .

1976 This model

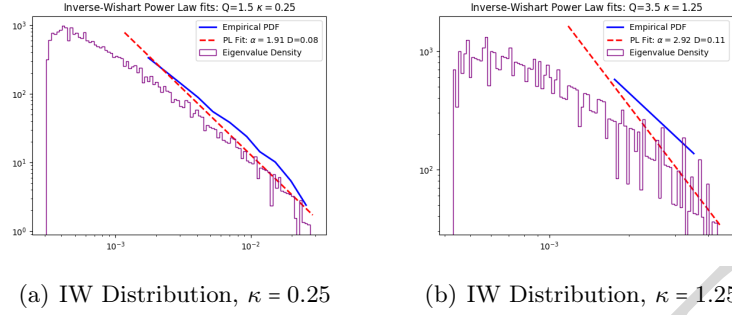


Figure 11: Example Inverse-Wishart (IW) distributions for  $\kappa = 0.25$  and  $\kappa = 1.25$ , along with Power Law (PL) fits of the generated distribution. Plots on Log-Log scale.

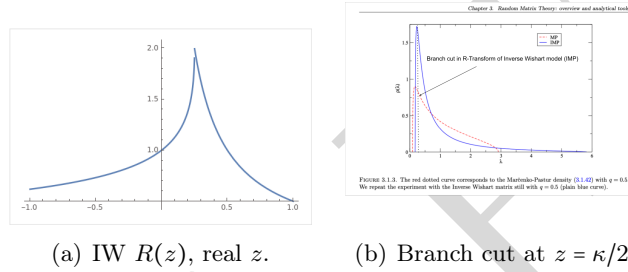


Figure 12: (a) The function  $R(z)$  of the Inverse Wishart model, with a singularity at  $z = \kappa/2$ . (b) The branch cut in the empirical spectral density, corresponding to the tail for  $\kappa = 0.5$ .

In Figure 11, we fit some typical layer ESDs to an IW distribution. When  $\kappa = 0.25$ , the fitted  $\alpha = 1.91$ , and the fit is a reasonably accurate model of the underlying Power Law distribution and the **WeightWatcher** PL fit. For  $\kappa = 1.25$ , the fitted  $\alpha = 2.92$  is larger, but the fit is not as good as a model. Generally speaking,  $\alpha$  scales with  $\kappa$ , but the free cumulants scale inversely with  $\kappa$ . So smaller  $\alpha$  will give larger free cumulants and therefore a larger  $\mathcal{Q}^2$ . Importantly, as seen in Figure 11(a), for  $\alpha \simeq 2.0$ , the IW model (with  $\kappa = 0.25$ ) is an effective simple model to illustrate the **SETOL** case of Ideal Learning.

Lets consider  $R(z)$  for the Inverse-Wishart model, denoted  $R(z)[IW]$ . To integrate this function, we require that it be analytic. At first glance, it may seem that that  $R(z)[IW]$  is not analytic because it has a pole at  $z = 0$  and because the square-root term  $\sqrt{\kappa(\kappa - 2z)}$  creates branch cut at and  $z = \kappa/2$  (and  $z = \infty$ ). Figure 12 presents this in two ways: Figure 12(a) shows the R-transform  $R(z)[IZ]$  for real  $z$ , highlighting its singular behavior and the location of the branch cut at  $z = \kappa/2$ ; and Figure 12(b) shows the corresponding branch cut in the ESD of the Inverse Wishart model (for  $\kappa = 0.5$ ). We select the branch cut starting at  $z = \kappa/2$  and ending at  $z = \infty$ , which allows us to at least formally defined the integral along the physically meaningful part of the ESD:

$$\mathcal{G}(\lambda)[IW] := \int_{\lambda_{min}^{ECS}}^{\lambda} R(z)[IW] dz, \quad (145)$$

noting that we expect  $\lambda_{min}^{ECS} \geq \kappa/2$ .

It turns out, however, that due to the branch cut in  $R(z)[IW]$ , the function  $\mathcal{G}(\lambda)[IW]$  is not analytic in the domain we need. To correct for this, we will instead model the Layer Quality-Squared using the modulus of  $\mathcal{G}(\lambda)[IW]$ ,

$$|\mathcal{G}(\lambda)[IW]| := \sqrt{\mathcal{G}(\lambda)[IW]^* \mathcal{G}(\lambda)[IW]} \quad (146)$$

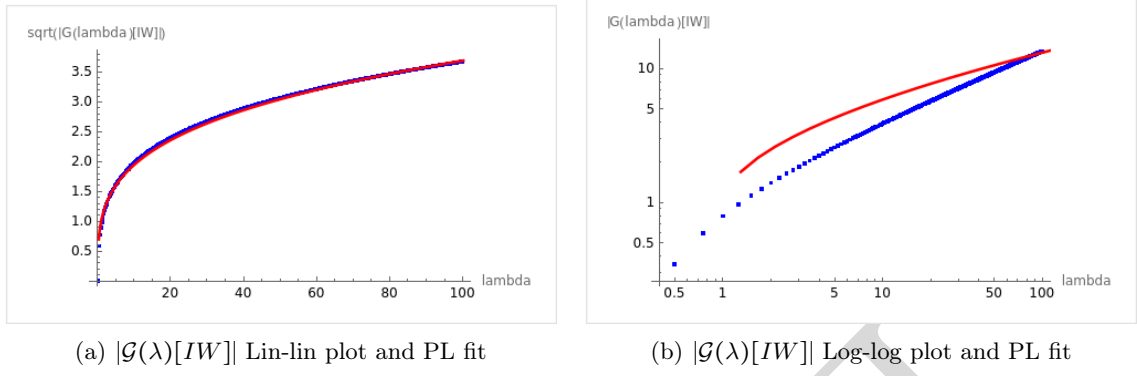


Figure 13: Behavior of  $\mathcal{G}(\lambda)[IW]$  for the Inverse Wishart (IW) model, with a Power Law (PL) fit (red),  $|\mathcal{G}(\lambda)[IW]| \approx 1.138\lambda^{0.539}$ . (a) Lin-lin plot. (b) Log-log plot.

where  $\mathcal{G}(\lambda)[IW]^*$  is the complex conjugate of  $\mathcal{G}(\lambda)[IW]$ . This is somewhat involved, so we present the full calculation in Appendix A.7

Figure 13 plots  $|\mathcal{G}(\lambda)[IW]|$  on Lin-lin and Log-log plots, on the range  $\lambda \in (0.25, 100)$ , and fits it to a PL. The fit follows the general trend of the function, but it is not terribly accurate. Still, from this plot, we can see that the Layer Quality-Squared has the same general trend as AlphaHat and/or a Shatten Norm.

#### 5.4.5 Levy-Wigner Models and the AlphaHat Metric

Here, we consider Levy-Wigner Models. We show how to obtain the WeightWatcher AlphaHat metric by modeling the near VHT cases with an approximation to a Levy distribution at  $\alpha \approx 2$ .

We do this because the AlphaHat metric has been developed to adjust for *Scale* anomalies that arise from issues like Correlation Traps, making Alpha smaller than expected. The Levy-Wigner (LW) model treats  $\mathbf{X}$  as if it were a Wigner matrix (and not actually a correlation matrix), and the  $\alpha$  is different but related to  $\alpha$  above in Table 5. The ESD follows a Levy-Stable distribution, where  $a$  is a shift parameter, and  $b$  is a complex phase factor depending on 2 real factors,  $\beta$  and  $\gamma$ . Strictly the ESD for an LW model,  $\rho_{LW}(\lambda)$ , is defined by its characteristic function (i.e., the Fourier Transform of  $\rho_{LW}(\lambda)$ ), but we can note that the ESD is VHT,  $\rho_{LW}(\lambda) \sim \lambda^{-\alpha-1}$ , and that when  $\alpha_l \approx 1$ , the ESD resembles a PL HT ESD with  $\alpha \approx 2$ .

For case of Ideal Learning, we choose to *model* the R-transform of our Fat-Tailed HT ESDs as

$$R(z)[HT] = bz^{\alpha-1}, \quad \alpha \approx 2 \quad (147)$$

where  $b$  is an unspecified constant (possibly negative and/or complex). Notice that when  $\alpha \approx 2$ , our model is close to the LW model,  $R(z)[HT] \approx R(z)[LW]$  (and gives a Cauchy distribution if we choose  $b = a - i\gamma$ ).

Integrating  $R(z)[HT]$ , and (as above) taking the approximation  $\lambda_{min}^{ECS} \sim 0$ , we obtain (formally)

$$\mathcal{G}(\lambda)[HT] = \frac{b}{\alpha} \lambda^\alpha. \quad (148)$$

If we now choose  $b = \alpha = 2$ , then  $\bar{Q}^2$  takes the form of a Shatten Norm (squared)

$$\bar{Q}^2 = \frac{1}{M} \sum_i \lambda^\alpha. \quad (149)$$

Taking the logarithm of  $\mathcal{G}(\lambda)[HT]$ , we obtain

$$\log \mathcal{G}(\lambda)[HT] = \log \frac{b}{\alpha} + \frac{\alpha}{\log} \lambda \quad (150)$$

2021 As with the Inverse-Wishart (IW) model, we can derive a formal expression for **AlphaHat** using  
 2022 the LW model. To do so, let us approximate  $\bar{Q}^2$  by the largest term in the sum over  $\mathcal{G}(\lambda)$ , and  
 2023 then let  $\lambda = \lambda_{max}$ , giving

$$\hat{\alpha} = \log_{10} \bar{Q}^2 \approx \alpha \log \lambda_{max}. \quad (151)$$

2024 We present this as a formal example, noting that is slightly different from the result for the IW  
 2025 model, Eqn. ???. We do not claim this is a valid empirical model, as we have not attempted to fit a  
 2026 real-world ESD to Levy-stable distribution. We leave this to a future study, noting, however, there  
 2027 has been some work doing such fits [?].

2028 Ideally, we would like to have an rigorous expression for  $R(z)$  not just in the case of Ideal  
 2029 Learning but also for the entire Fat-Tailed Universality class. This is non-trivial to obtain and we  
 2030 will attempt this in a future work. Fow now, we will take a different approach, and evaluate  $R(z)$   
 2031 explicitly using numerical techniques.

## 6 Empirical Studies

In this section, we present empirical results. Our goals are to justify key technical claims, including key assumptions underlying our SETOL approach, and to illustrate the behavior of SETOL with respect to various parameters and hyperparameters. Importantly, it is *not* our goal to demonstrate that layer PL exponent Alpha and AlphaHat perform well for diagnostics and predicting model quality for SOTA NN models, as that has been demonstrated previously [?, ?, ?].

Since the SETOL theory presented in Section 5 is (effectively) a single layer theory, in order to carefully test (as opposed to simply use) SETOL, we need to limit the degree of inter-layer interactions present in the model. To do so, we consider a three-layer Multi-Layer Perceptron (MLP3), trained on MNIST [?]. We refer to the hidden layers as “FC1 and “FC2. Their output sizes and parameter counts are shown in Table 6.

Layer	Units	Weight Parameters	% of total
FC1	300	$768 \times 300 = 230,700$	88.2%
FC2	100	$300 \times 100 = 30,000$	11.4%
out	10	$10 \times 100 = 1000$	0.38%

Table 6: Dimensions of each FC layer in the MLP3 model, along with weight matrix parameter count and fraction of the total.

The following are the main topics we consider.

**6.1 Model Quality: HTSR phenomenology.** The HTSR phenomenology provides a metric of model quality in the form of the PL exponent  $\alpha$ .<sup>36</sup> In particular, smaller values of  $\alpha$  (e.g., values of  $\alpha$  closer to 2 than 3 or 4) should correspond to better models, e.g., having smaller test errors; and a minimal error should be obtained when  $\alpha = 2$ . See Section 6.1.

**6.2 Effective Correlation Space.** The SETOL theory is based on the notion of an Effective Correlation Space, in which the learning and generalization occurs. This is the low-rank subspace  $\mathbf{W}^{ecs}$  of each layer  $\mathbf{W}$  that approximates the teacher  $\mathbf{T}$ . In particular, our measure of model quality should be restricted to  $\mathbf{W}^{ecs}$ . The Effective Correlation Space can be identified from the tail of the ESD,  $\rho_{tail}(\lambda)$ , and it can be chosen according to one of several related Model Selection Rules. See Section 6.2.

**6.3 Evaluating the Trace-Log Condition.** In the HTSR phenomenology, when a model is very well-trained, the layer PL exponent  $\alpha \simeq 2$ . In the SETOL theory, when a model is very well-trained, the eigenvalues in the tail will satisfy the Empirical Trace-Log Condition, given in Eqn. (125). In Section 6.3, we provide a detailed analysis of this effect.

**6.4 Correlation Traps.** Recall, from Section 3.3.1, that if a layer weight matrix  $\mathbf{W}$  has a Correlation Trap (and, in particular, those arising from SGD training with very large learning rates) then it is likely that the test (and train) accuracy will be degraded, and  $\alpha$  will drop below its optimal value. See Section 6.4 for an empirical demonstration of this.

**6.5 Overloading and Hysteresis Effects.** The experiments described so far validate that SETOL makes valid predictions in the  $\alpha \gtrsim 2$  range. Beyond that point, SETOL only predicts

<sup>36</sup>Prior work has shown that the AlphaHat metric ( $\hat{\alpha}$ ) accurately describes variations in model quality as a function of architecture changes [?]. Since we do not vary the depth of the model in our evaluations, the Alpha metric ( $\alpha$ ) is of interest in this setting.



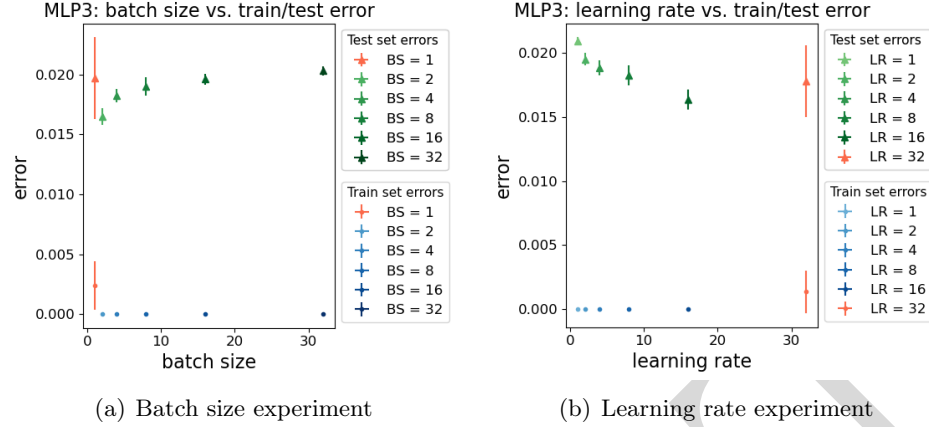


Figure 14: Train / test errors in the MLP3 model as a function of batch size, and learning rate. Observe the inverse relationship between batch size (a) and learning rate (b). As batch size decreases, test error generally decreases, until batch size reached  $bs = 1$ . Similarly, as learning rate increases, test error decreases until  $lr = 32\times$  the SGD default value of 0.01.

“atypicality, in the sense of spin glass theory [?]. See Section 6.5 for an examination of how the MLP3 behaves when it is pushed further out of that range of validity, e.g., by training only one layer, while keeping the others frozen. In particular, we compare results when a single layer is either over- vs under-parameterized.

We trained the MLP3 model independently, using both the Tensorflow 2.0 framework (using the Keras api, and on Google Colab) and pytorch, with the goal of consistent, reproducible results. Each setting of batch size, learning rate, and trainable layer was trained with 5 separate starting random seeds, and error bars shown in plots below represent one standard deviation taken over these 5 random seeds. Each training run used the same early stopping criteria on the train loss: training was halted when train loss did not decrease by more than  $\Delta E_{train} = 0.0001$ , over a period of 3 epochs. In doing so, each model was trained with a different number of epochs; and, at the end, the best weights were chosen for the model. See Appendix A.5 for more details on the MLP3 model and the training setup. We provide a Google Colab notebook with the exact details, allowing the reader to reproduce the results as desired.

The dominant generalizing components of  $\mathbf{W}$  reside in  $\mathbf{W}^{ecs}$  such that it captures the functional contribution of  $\mathbf{W}$  to the NN; and thus

## 6.1 HTSR Phenomenology: Predicting Model Quality via the Alpha metric

Here, we want to determine how the quality of our MLP3 model varies with the Alpha metric. From previous work [?, ?, ?], we expect that Alpha metrics for the FC1 and FC2 layers should be well-correlated with the test accuracy, while varying some suitable training knob, such as learning rate or batch size, that can modulate the test accuracy.<sup>37</sup>

We vary the batch size from small to large, i.e.,  $bs \in [1, 2, 4, 8, 16, 32]$ , following the setup of previous work on the HTSR phenomenology [?]. We expect similar effects by varying the learning rate, as it is known that small batch sizes correspond directly to large learning rates [?, ?]. Thus, we conducted a second set of experiments where the learning rate was varied by a factor of

<sup>37</sup>Since we do not change the depth of the model here, we expect the Alpha metric to follow the AlphaHat metric, also predicting the test accuracies [?].

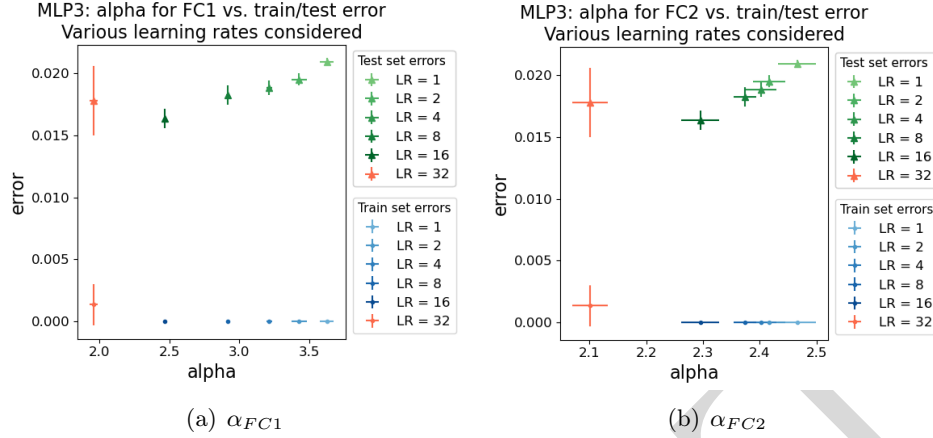


Figure 15: Train / test errors in the MLP3 model in the *Learning Rate* experiment as a function of  $\alpha_{FC1}$  (a) and  $\alpha_{FC2}$  (b). Observe the regular downward progression of  $\alpha$  and error as the learning rate increased in both (a) and (b). When learning rate was 32 $\times$ , (shown in red),  $\alpha_{FC1}$  fell below 2, coinciding with a drastic increase in both train and test error. The results here almost perfectly replicate those of the Batch Size experiment, shown in Figure 16.

[1 $\times$ , 2 $\times$ , 4 $\times$ , 8 $\times$ , 16 $\times$ , 32 $\times$ ], relative to the SGD default value of 0.01. Adjusting the learning rate or batch size allows us to systematically vary the layer PL exponent  $\alpha$  between roughly 2 and 4, i.e., within the range in which SETOL should make the most reliable predictions. As an added benefit, it also allows us to use the very small batch size of 1 to force the model into a state of over-regularization, which we also analyze below.

Consider Figure 14, which plots the final train and test accuracies as a function of the hyperparameter (batch size or learning rate) used during training for the MLP3 model. Figure 14(a) varies batch size, and Figure 14(b) varies learning rate. Recall that error bars represent one standard deviation taken over 5 independent starting random seeds. In Figure 14(a), we see that by decreasing the batch size ( $bs$ ), and holding other knobs constant, we can systematically improve the train and test accuracy, up to a point. In particular, for  $bs \geq 2$ , both the test and train accuracies increase with decreasing batch size, consistent with previous work [?]. Further decrease beyond  $bs = 2$  leads to *lower* model quality, i.e., larger error and larger error variability. In Figure 14, we see that increasing the learning rate ( $lr$ ) by a factor  $x$  has an exactly analogous effect as decreasing the batch size by  $1/x$ .<sup>38</sup>

The transition between  $lr = 16\times$  (or  $bs = 2$ ), which is locally optimal for the setting of other hyperparameters, and  $lr = 32\times$  normal (or  $bs = 1$ ), which is not, provides a demonstration of a distinct change in the behavior of **Alpha**, concordant with the sudden increase in the error and error variability. To explore this in the context of SETOL, consider Figure 15 and Figure 16, which plot error as a function of **Alpha**, for different learning rates and batch sizes, respectively.

Figure 15 plots the **Alpha** metrics  $\alpha_{FC1}$  and  $\alpha_{FC2}$ , as learning rate is varied, demonstrating that both metrics are well-correlated with the test accuracies, for all learning rates less than 16 $\times$  normal. In particular, as we drive **Alpha** in FC1 down to an Ideal value of  $\alpha \simeq 2$ , the test error decreases monotonically (Figure 15(a)). Beyond that point, further decrease of the batch size sees **Alpha** decrease below its Ideal value of 2 in the FC1 layer. This corresponds not only with *higher* errors, but also with *larger* error bars, on both train and test error. The dramatic increase in train error and train error variability is particularly telling, because it suggests that when  $\alpha_{FC1}$  passes

<sup>38</sup>One could, of course, mitigate this by fiddling with other knobs of the training process, but that is not our goal. Our goal here is not to use a toy model to demonstrate the properties and predictions of SETOL.

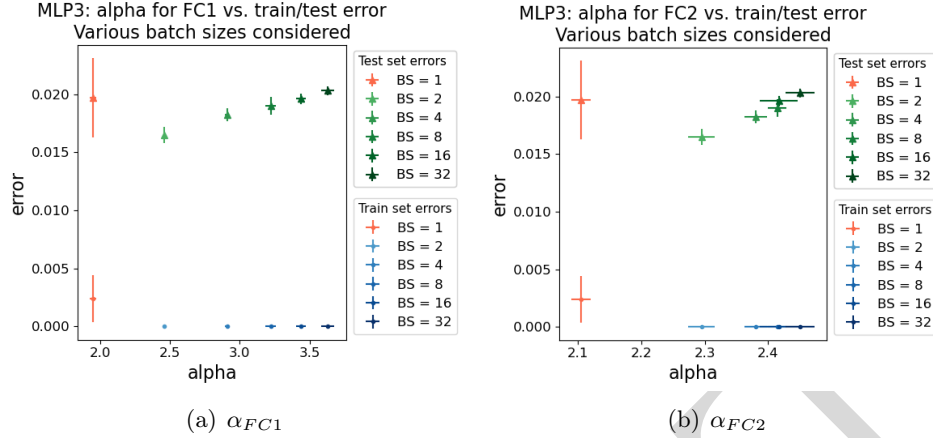


Figure 16: Train / test errors in the MLP3 model in the **Batch Size** experiment as a function of  $\alpha_{FC1}$  (a) and  $\alpha_{FC2}$  (b). Observe the regular downward progression of  $\alpha$  and error as the batch size decreased in both (a) and (b). When batch size was 1, (shown in red),  $\alpha_{FC1}$  fell below 2, coinciding with a drastic increase in both train and test error. The results here almost perfectly replicate those of the Learning Rate experiment, shown in Figure 15.

below 2, the model enters into a “glassy state, and is unable to relax down to 0 train error.

In Figure 15(b), we consider **Alpha** for FC2, and we see that  $\alpha_{FC2}$  approaches 2, but does not reach it, even for  $lr = 32\times$ . This failure to achieve  $\alpha_{FC2} \approx 2$ , along with the much greater size of FC1, (See Table 6,) suggests that FC1 is the critical layer for the models performance. This also highlights some of the interplay between the layers, (which is not captured by a single layer theory) – as  $\alpha_{FC1}$  has narrow error bars throughout,  $\alpha_{FC2}$  shows much more variation by way of its wider error bars. Thus, while model accuracy kept improving as learning rate increased up to  $16\times$ , this was likely driven by a better  $\alpha_{FC1}$ , more than by  $\alpha_{FC2}$ .

In Figure 16, we consider batch size, and we see a near identical replication of these results, in terms of the relation of train error, test error and **Alpha** in the two layers. Consequently, the remainder of the experiments will focus on the learning rate experiment, as both produced substantially the same results.

## 6.2 Testing the Effective Correlation Space

Here, we will address the question:

How shall we *test the assumption* of the Effective Correlation Space?

Recall that the SETOL theory estimates model quality by evaluating the ST Generalization Error as an integral over the theoretical training data  $\xi$ . This integral assumes each layer weight matrix can be replaced with an effectively lower rank form, i.e.,  $\mathbf{W} \rightarrow \mathbf{W}^{ecs}$ , corresponding to the span of the eigencomponents defined by the tail of the ESD,  $\rho_{tail}(\lambda)$ . In the HTSR phenomenology, the tail is defined by the fact that  $\rho_{tail}(\lambda)$  follows a PL distribution, above some minimal value  $\lambda_{min}$ . In our SETOL theory, the tail is defined by choosing the minimal value  $\lambda_{min}$  to satisfy the Empirical Trace-Log condition. These methods of realizing  $\mathbf{W}^{ecs}$  are essentially *Model Selection Rules* (MSRs) for the Effective Correlation Space. Importantly, in neither approach is  $\lambda_{min}$  just some “rank parameter to be chosen by yet some other MSR on the basis of rank, or magnitude alone, (that, in particular, does not know about HTSR or SETOL, which consider the *shape* of the ESD).

Thus, to test the assumption of the Effective Correlation Space, we want to show that the models predictions are in fact controlled predominantly by the tail, where the specific choice of the rank parameter depends on HTSR or SETOL as we expect. We can emulate this theoretical construct and estimate (trends in the) test accuracies by evaluating the train and/or test accuracies of the trained MLP3 model – after replacing the MLP3 layer weight matrices  $\mathbf{W}_{FC1}$  and  $\mathbf{W}_{FC2}$  with a low-rank approximation consisting of *only* the tail:

$$\begin{aligned}\mathbf{W}_{FC1}^{ecs} &:= P_{tail} \mathbf{W}_{FC1} \\ \mathbf{W}_{FC2}^{ecs} &:= P_{tail} \mathbf{W}_{FC2},\end{aligned}$$

where  $P_{tail}$  is a projection operator selecting only the tail of the ESD with TruncatedSVD. (That is, we will use the low-rank TruncatedSVD approximation at the inference step, not at the training step, as is more common.) A Truncated model is one whose weight matrices  $\mathbf{W}_*$  have been replaced by truncated matrices  $\mathbf{W}_*^{ecs}$ . We denote the difference between the original models accuracy and the Truncated models train and test accuracy as  $\Delta E_{train}$  and  $\Delta E_{test}$ , respectively:

$$\begin{aligned}\Delta E_{train} &:= E_{train}(\mathcal{D}) - E_{train}^{ecs}(\mathcal{D}) \\ \Delta E_{test} &:= E_{test}(\mathcal{D}) - E_{test}^{ecs}(\mathcal{D}),\end{aligned}$$

where  $E_{train}^{ecs}$  denotes the error of the TruncatedSVD model on the training portion of the dataset  $\mathcal{D}$ , and  $E_{test}^{ecs}$  denotes corresponding test error for the TruncatedSVD model.

**The PowerLaw and TRACE-LOG Model Selection Rules** If we use good MSRs, then we expect that  $\Delta E_{train} \rightarrow 0$  and  $\Delta E_{test} \rightarrow 0$  as the models approach Ideal Learning. We consider the following MSRs,<sup>39</sup> which are associated with the HTSR and SETOL approaches.

- The **PowerLaw** MSR: All eigenvalues lying in the tail of the ESD,  $\lambda_i \geq \lambda_{min}^{PL}$ , where  $\lambda_{min}^{PL}$  is the start of the PL tail, as determined by the **WeightWatcher** PL fit, which is based on [?].
- The **TRACE-LOG** MSR: All eigenvalues lying in the tail of the ESD, such that they satisfy the Trace-Log Condition, i.e.,  $\lambda_i \geq \lambda_{min}^{|detX|=1}$ , where  $\prod_{i: \lambda_i \geq \lambda_{min}^{|detX|=1}} \lambda_i \simeq 1$ .

### 6.2.1 Train and test errors by epochs

To see how the Effective Correlation Space forms, we plot how  $\Delta E_{train}$  and  $\Delta E_{test}$  evolve over training, for each of the various learning rates considered.<sup>40</sup>

We start with the effect of the **PowerLaw** MSR. See Figure 17, where we see that  $\Delta E_{train}$  and  $\Delta E_{test}$  generally trend downwards as they approach minimum train error. When the learning rate is larger, the models converge more quickly, and  $\Delta E_{train}$  and  $\Delta E_{test}$  also converge to lower values. Recall from Figure 14(b) that  $lr = 16\times$  had the lowest test error. In Figure 17(e), we see that it also has the lowest  $\Delta E_{train}$  and  $\Delta E_{test}$ . A lower  $\Delta E_{train}$  or  $\Delta E_{test}$  means that more of the models correct predictions are due to the low rank tail, meaning that the tail generalizes better, and we see here that when the tail generalizes best, the model was the most accurate.

In each plot, we also see that the error bars are wide early on, before suddenly becoming much narrower. This transition is more visible in the larger learning rates shown in 17(d)–17(f), but can also be seen in 17(a)–17(c), albeit less clearly. Most interestingly of all, this transition is preceded

<sup>39</sup>We considered other MSRs that do not “know about HTSR or SETOL, but they (expectedly) perform in trivial or uninteresting ways for testing the assumption of the Effective Correlation Space. Thus, we are not introducing just some arbitrary low-rank approximation, as is common, but instead that the specific SETOL-based MSR matters.

<sup>40</sup>When batch size was varied, results did not significantly differ, and so they are omitted.

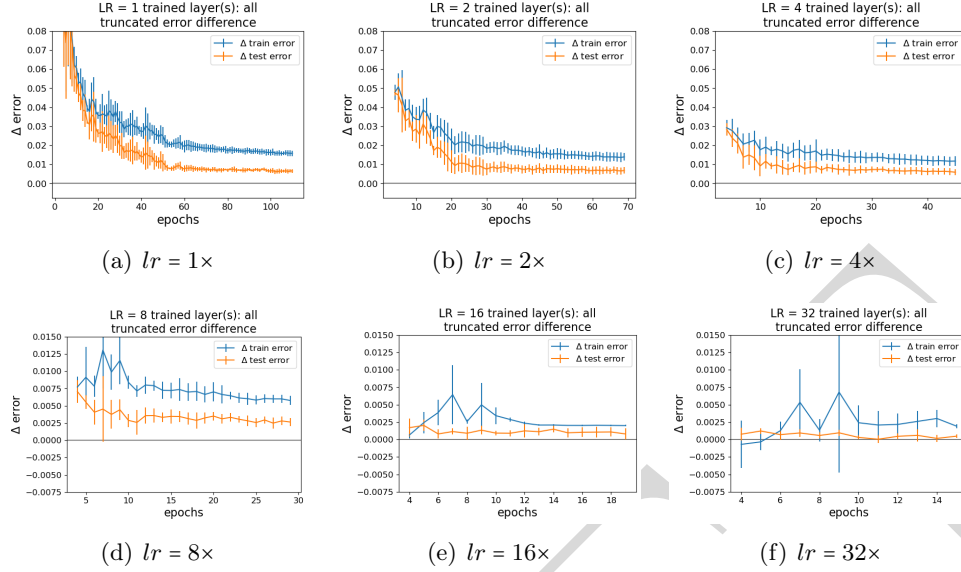


Figure 17:  $\Delta E_{train}$  (blue) and  $\Delta E_{test}$  (orange) for various learning rates, using the **PowerLaw** MSR. As learning rate increases we can see that  $\Delta E_{train}$  and  $\Delta E_{test}$  both tend towards lower asymptotic minima, which they reach after fewer epochs of training. We can also see that (after the first few epochs,)  $\Delta E_{train}$  (blue) is always higher than  $\Delta E_{test}$  (orange). Observe that in the bottom row (d–f) the yaxis is contracted to make the variation more visible. In (f) we can see that as learning rate surpasses its optimal setting, the gap between  $\Delta E_{train}$  and  $\Delta E_{test}$  begins to increase again, and has wider error bars, suggesting that the excessively large learning rate is disrupting the MLP3s ability to learn the Effective Correlation Space.

by a brief period, sometimes a single epoch, in which the error bars are drastically wider, in a way that is reminiscent of a first-order phase transition. Again, this phenomenon can be seen most clearly in 17(d)–17(f).

We next consider the effect of the **TRACE-LOG** MSR. See Figure 18, which also shows the development of  $\Delta E_{train}$  and  $\Delta E_{test}$  over epochs, where we see a very different pattern in the train error and test error. The difference in the train error is because, as the model is untrained in the early epochs, the **TRACE-LOG** MSR *over-estimates* the tail by choosing a  $\lambda_{min}$  that is too small. Thus,  $\Delta E_{train}$  actually increases to its asymptotic value at the final epoch. In the earliest epochs, the truncated train error is even *less* than the full MLP3 models error, suggesting that signal is forming in the large eigenvalues in these early epochs, but is swamped by the randomness of the early initial weights, some of which is then removed by truncation. As epochs progress, this effect disappears. Here again we can see the “phase-transition-like behavior of the train error, as the error bars are wide early on, up to a transition having an abnormally large error bar, after which they stabilize.

Perhaps most interestingly of all, we see that under the **TRACE-LOG** MSR,  $\Delta E_{test}$  is flat *throughout training*, and for all learning rates. Considering that this implies that there is *no* point in training where the **TRACE-LOG** tail generalizes badly, this is a rather striking observation. It also bears noticing that under the **PowerLaw** MSR, (Figure 17,) this is decidedly not the case, meaning that a small, but significant amount of generalization comes from the gap between these tails, but none at all comes from the bulk. [confusing]

There are two final points of comparison between Figures 17 and 18. First, although it appears in Figure 18 that  $\Delta E_{train}$  converges to a larger value than in Figure 17, this is because the scale



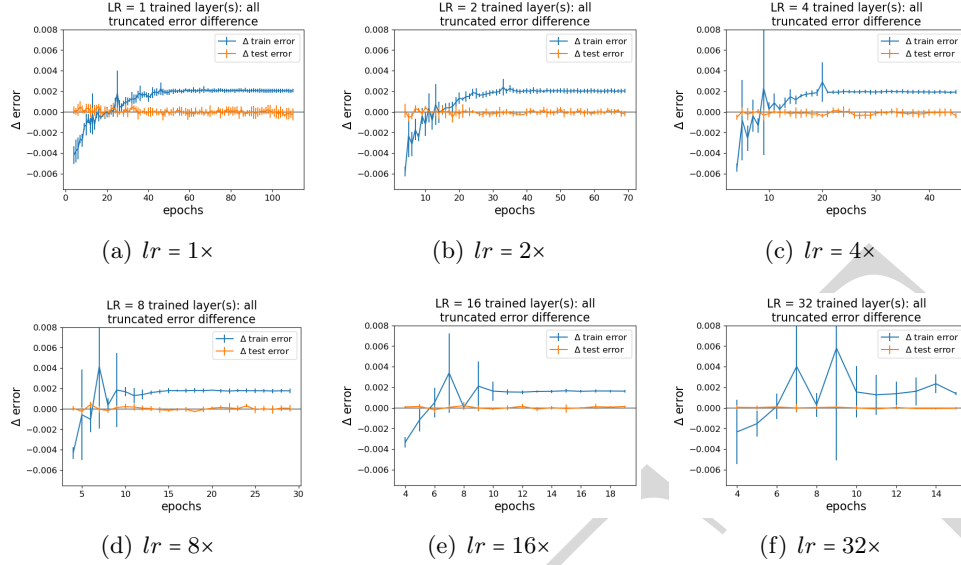


Figure 18:  $\Delta E_{train}$  (blue) and  $\Delta E_{test}$  (orange) for selected learning rates, using the TRACE-LOG MSR. For all learning rates,  $\Delta E_{test}$  is centered on 0, meaning that the TRACE-LOG Effective Correlation Space explains almost all variation in out-of-sample predictions, but it does *not* explain all of the training set predictions (blue). NOTE: The y axis is the same in all plots, and is much narrower than in Figure 17. In (a–e) we can see that  $\Delta E_{train}$  converges to approximately 0.002. Compare with Figure 17(e), which reaches a minimum of approximately 0.0025. As in Figure 17(f), the learning rate of 32 $\times$  normal disrupts the MLP3s ability to discover the TRACE-LOG Effective Correlation Space.

of the y-axis is 10 $\times$  smaller. That is, the PowerLaw MSR is biased towards over-estimating  $\lambda_{min}$ , which means it over-truncates, producing a larger  $\Delta E_{train}$  or  $\Delta E_{test}$  than the TRACE-LOG MSR, which is biased towards under-estimating  $\lambda_{min}$ . Second, in both Figure 17 and 18, we can see that  $\Delta E_{test}$  is consistently lower than  $\Delta E_{train}$ . Clearly, truncating has a larger effect on train predictions, meaning that no matter how long the model is trained, some of the train predictions are still derived from the bulk. Yet, the test predictions are far less affected, meaning that the Effective Correlation Space contributes to the models ability to generalize.

## 6.2.2 Truncation and Generalization

Given that  $\Delta E_{test}$  is always lower than  $\Delta E_{train}$  for the PowerLaw MSR, and similarly for TRACE-LOG after a certain point in training, it is clear that the Effective Correlation Space has something to do with generalization. However, this leaves open the question of what role precisely Alpha plays. In Figure 19, we plot  $\Delta E_{train}$  and  $\Delta E_{test}$  with the PowerLaw MSR as a function of Alpha (rather than epochs) for layers FC1 and FC2, as well as the Generalization Gap—that is,  $E_{test} - E_{train}$ . (Recall Eqn. 83.) Learning rate is not explicitly shown, but its effect can be seen in the clusters of points that each learning rate generates.

In both layers,  $\Delta E_{train}$  and  $\Delta E_{test}$  steadily decrease with  $\alpha_{FC1}$ , until it passes below 2, after which the relation deteriorates somewhat. This is especially prominent in FC2. Recall from Figure 15, Section 6.1, that when  $\alpha_{FC1}$  passed below 2, the train error and test error both increased and exhibited larger variability. From this we interpret Alpha as a measure of *regularization* (which is consistent with its introduction as a measure of implicit self-regularization [?]). Regularization has the effect of keeping train and test accuracy closer together, and generally, as Alpha in the



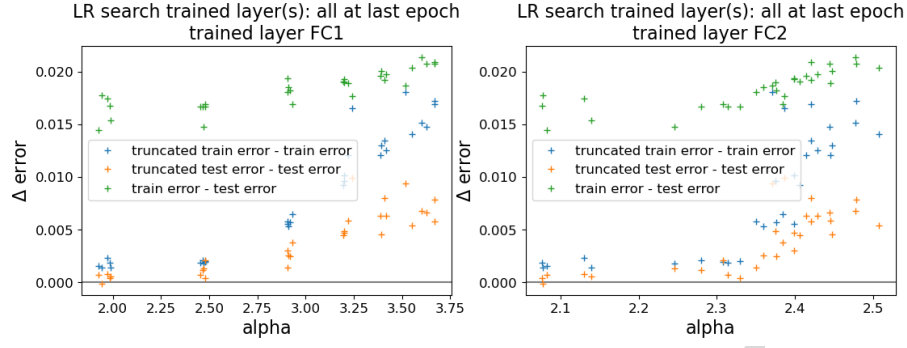


Figure 19: Train and test error gaps using the **PowerLaw** MSR, as a function of  $\alpha$  in the FC1 and FC2 layers of MLP3 models, at the *final epoch* of training. We can see that as  $\alpha$  decreases towards 2, (right to left),  $\Delta E_{train}$  and  $\Delta E_{test}$  generally decrease as well, meaning that the closer  $\alpha$  is to 2, the more the Effective Correlation Space explains the train and test predictions. The gap between un-truncated train and test error, (Eqn.83,) generally decreases as well until  $\alpha_{FC1} < 2$ .

dominant layer decreases towards 2 from above, the train-test error gap decreases.

### 6.3 Evaluating the Trace-Log Condition

Having established that the PL tail of the ESD, defined by eigenvalues above  $\lambda_{min}^{PL}$ , is a major factor in determining model quality in the MLP3 model, we now examine how well the Trace-Log Condition compares with it. In particular, we demonstrate that when the tail of a layer ESD is described well by the HTSR phenomenology, i.e., when it is well-fit by a PL with  $\rho(\lambda)_{tail} \sim \lambda^{-\alpha}$ , with PL exponent  $\alpha \simeq 2$ , then the eigenvalues in the tail defined by the PL fit, i.e.,  $\lambda \geq \lambda_{min}^{PL}$ , also satisfy the Trace-Log Condition of Eqn. 125—a key assumption of the **SETOL** theory. This is a rather remarkable empirical result that couples HTSR and SETOL; it has its basis in our SETOL derivation; and it provides the basis for an inductive principle that is based on the product of eigenvalues rather than an eigenvalue gap.

We can denote the eigenvalue that best fits the Trace-Log Condition as  $\lambda_{min}^{|detX|=1}$ . Then, to measure how well this condition holds, we can compute

$$\Delta\lambda_{min} = \lambda_{min}^{PL} - \lambda_{min}^{|detX|=1}. \quad (152)$$

In Sections 6.3.1 and 6.3.2, we will see the trend that as  $\alpha$  approaches 2,  $\lambda_{min}^{PL}$  and  $\lambda_{min}^{|detX|=1}$  also approach one another, and hence  $\Delta\lambda_{min}$  goes to 0, from above, both for our toy MLP3 model as for SOTA models. In our MLP3 model, we will see that a crossing of the equality condition coincides with over-regularization and a degradation in model accuracy. In pre-trained ResNet[?], VGG[?] and ViT[?] models, we will also see, empirically, that in general  $\Delta\lambda_{min}$  remains positive, just as  $\alpha$  remains above 2.

#### 6.3.1 The MLP3 model

Consider Figure 20, which shows  $\lambda_{min}^{PL}$  and  $\lambda_{min}^{|detX|=1}$  in the FC1 layer of three MLP3 models, each sharing a common starting random seed, that were trained with the largest learning rates. The  $\lambda_{min}^{PL}$  and  $\lambda_{min}^{|detX|=1}$  eigenvalues are marked by red and purple vertical lines, respectively; and thus  $\Delta\lambda_{min}$  is the distance between red and purple lines. As learning rate increases, the red and purple

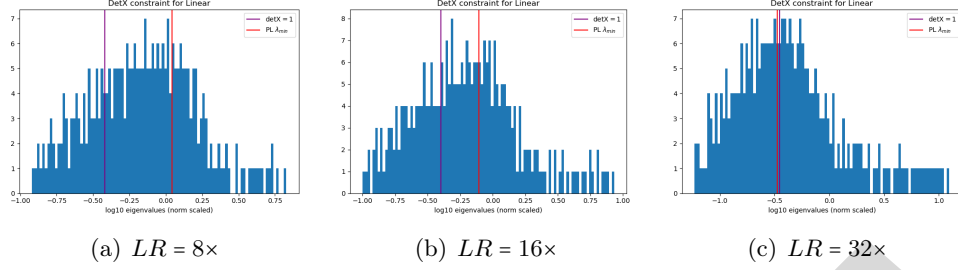


Figure 20: Log-Linear ESDs for three learning rates in the FC1 layer of MLP3. The red line shows  $\lambda_{min}^{PL}$ , and the purple line shows  $\lambda_{min}^{|detX|=1}$ . Observe that the purple line is to the left of the red line, but as the LR increases they move closer together. However, when LR is  $32\times$ , where both train and test accuracy suffered (c), the red line is to the left of the purple line. This is often a signature of an Over-Regularized layer, and indeed the FC1 layer in this model had  $\alpha < 2$ . (See Figure 15(a) in Section 6.1.)

lines draw closer, and they are closest for  $lr = 32\times$  (Figure 20(c)). (Compare this with Figure 15, Section 6.1, which shows that this corresponds with an increase in test accuracy, up to  $lr = 16\times$ , but at  $lr = 32\times$  Alpha fell below 2 and accuracy suffered.) In Figure 20(a)-20(b), the purple line is left of the red line; but in Figure 20(c), the red and purple lines cross, such that  $\lambda_{min}^{PL} < \lambda_{min}^{|detX|=1}$ . This is analogous to the case where  $\alpha$  crosses below 2. This suggests that the absolute Trace-Log is minimized when  $\alpha \simeq 2$ , which (remarkably) is exactly when the HTSR phenomenology predicts the layer is Ideal.

(Observe that Ideal does not necessarily mean optimal under a finite sized training set, but rather that the finite sized system behaves the most similarly to its infinite limit.)

We can compare  $\alpha$  and  $\Delta\lambda_{min}$  more broadly by plotting  $\Delta\lambda_{min}$  directly as a function of  $\alpha$  in a single plot spanning all random seeds and learning rates or batch sizes. This is shown in Figure 21. Critical values of  $\alpha = 2$  and  $\Delta\lambda_{min} = 0$  are shown as vertical and horizontal red lines, respectively. Values for various learning rates are plotted for layer FC1 (Figure 21(a)) and FC2 (Figure 21(b)), as well as for various batch sizes in layer FC1 (Figure 21(c)) and FC2 (Figure 21(d)).

For layer FC1 (Figures 21(a) and 21(c)), in both cases we see near-linear march towards the critical tuple of  $(\alpha, \Delta\lambda_{min}) = (2, 0)$ . In addition, passing this critical value coincides with diminished train and test accuracy, (recall Figure 14), suggesting that just as  $\alpha = 2$  is a threshold of over-regularization,  $\lambda_{min}^{PL} < \lambda_{min}^{|detX|=1}$  may be as well. Since FC1 is the dominant layer, comprising roughly 8/9 of the weights of the model, (Table 6,) we expect FC1 to most closely match the performance of the model as a whole.

For layer FC2, which comprises roughly the other 1/9 of the models weights, there is a similar coevolution, but it is weaker. As learning rate or batch size exceeds their critical values, rather than going to  $(2, 0)$  as in FC1, we instead see that the relationship simply breaks down, with the gap growing larger even as  $\alpha$  decreases. Given that FC1 has passed the critical  $\alpha = 2$  threshold, we conjecture that the breakdown of the relationship between  $\alpha$  and  $\Delta\lambda_{min}$  is due to FC1 becoming atypical.

### 6.3.2 State-of-the-Art (SOTA) models

Here, we consider SOTA models, in particular VGG pre-trained models [?], the ResNet series [?], the ViT series [?], and the DenseNet series [?]. We show that as  $\alpha$  approaches 2, the Log-Trace Condition holds better and better, i.e.,  $\Delta\lambda_{min}$  approaches 0.

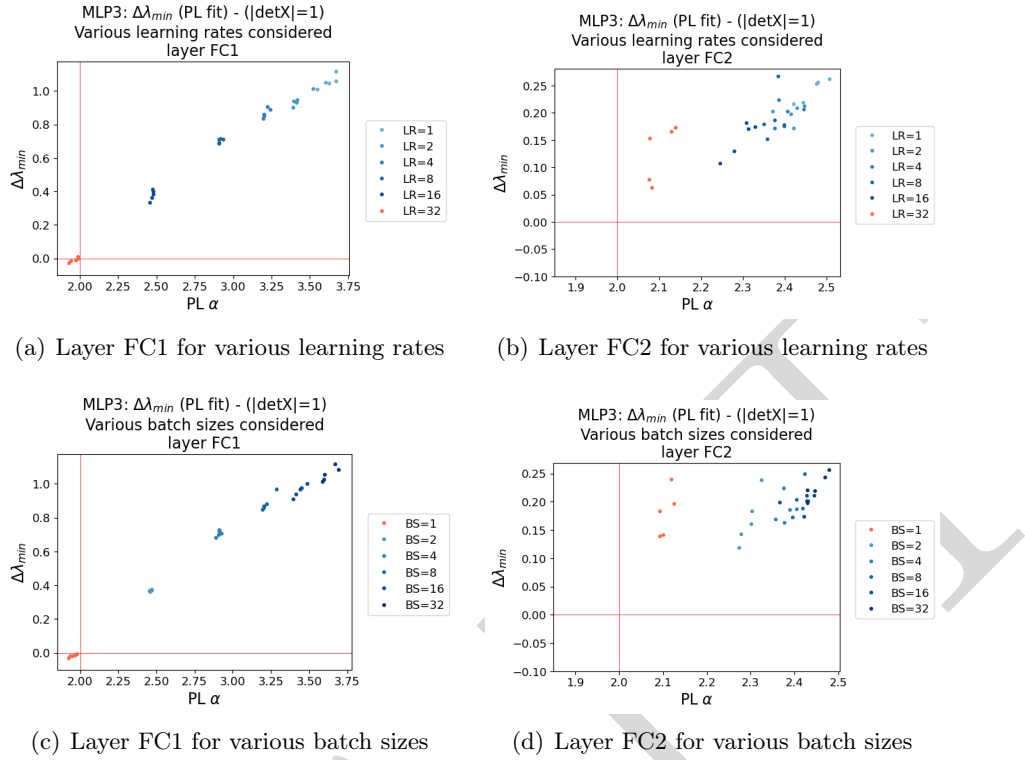


Figure 21: MLP3 Model: Comparison of the PL Alpha (x-axis), with the difference between  $\lambda_{min}^{PL}$  and  $\lambda_{min}^{|\det X|=1}$  (y-axis). The thin red lines indicate critical values of  $\alpha = 2$  and  $\Delta\lambda_{min} = 0$ . As learning rate increases (a–b) or batch size decreases (c–d), we can see that in layer FC1, which dominates the model, (See Table 6,)  $\alpha$  goes to 2, and  $\Delta\lambda_{min}$  goes to 0. Observe that both critical values are crossed at the most extreme hyper-parameter selection, (red,) corresponding with over-training. Layer FC2 shows a weaker tendency towards the critical values (b, d), and is disrupted at the most extreme hyper-parameter values (red).

Figure 22 plots  $\alpha$  versus the difference  $\Delta\lambda_{min}$ , (Eqn. 152). Layer matrices from all models in each series are pooled to generate the plots.<sup>41</sup> Notice that in Figure 22(a) – 22(d),  $\Delta\lambda_{min}$  approaches zero as  $\alpha \rightarrow 2$  from above. Individual points may have a large  $\Delta\lambda_{min}$  for an  $\alpha$  near to 2, but the overall trend is apparent. The rapid decrease of  $\Delta\lambda_{min}$  as  $\alpha$  approaches 2 from above, implies that the PL tail rapidly takes on a unit Trace-Log (if it doesn't already have it.)

As we saw in comparison of Figures 17 and 18, (Section 6.2.1,) the TRACE-LOG tail is generally larger, and always has highly generalizing components. Thus, it is plausible that as layers reach the limit of the amount of information that can be encoded in them, i.e. as  $\alpha$  goes to 2, the PowerLaw tail expands to fill the TRACE-LOG tail. This effect can be seen clearly in Figure 22 in the condensing of the “funnel shape.”

Recall from Figure 21 that layer FC1 dominated the model, (Table 6,) producing a clear progression of  $\alpha$  and  $\Delta\lambda_{min}$  towards (2,0), as a function of learning rate or batch size, whereas FC2 showed a slightly less clear relationship. In larger models having dozens of layers, we would not expect any one layer to dominate as thoroughly. Moreover, it is the architecture that varies

<sup>41</sup>For convolutional layers, the WeightWatcher tool first computes eigenvalues for all channels-to-channels linear operators separately, and then pools them in order to compute  $\alpha$ ,  $\lambda_{min}^{PL}$  and  $\lambda_{min}^{|\det X|=1}$ . For instance, in a  $64 \times 64 \times 3 \times 3$  weight tensor, there would be 9 separate linear operators of  $64 \times 64$ , giving 576 eigenvalues, which would then be pooled to compute  $\alpha$ ,  $\lambda_{min}^{PL}$  and  $\lambda_{min}^{|\det X|=1}$ .

between models in each series, not (necessarily) the hyperparameters, meaning that there would not be a straight line, as in Figures 21(a) and 21(c). However, with all of the layers contributing to varying degrees, we nevertheless see a clear trend in both plots of Figure 22. These results show how the single-layer SETOL theory extends from the MLP3 model, which is dominated by a single layer, to larger models where many layers interact in complex ways.

The overall pattern of relationship between  $\Delta\lambda_{min}$  and  $\alpha$  can also be seen in Figure 23, which shows plots for Large Language Models (LLMs) of the Falcon [?] and LLAMA [?] model families, for different numbers of parameters. Observe that each subfigure 23(a)–23(d) shows a single model, rather than a collection of models in a family, as in Figure 22. The y-axis is the same between models in the same family. As in Figure 22, there is a general outline of a “funnel shape” pointing towards the critical point (2, 0), with the exception that it is only reached in the case of LLAMA-65b, (Figure 23(d)). This suggests that these LLMs are larger than they necessarily need to be, consistent with prior work [?], but also that they are well guarded against Over-Regularized layers beyond the critical point ( $\alpha = 2$  and  $\Delta\lambda_{min} = 0$ ).

## 6.4 Inducing a Correlation Trap

Previous work on the HTSR phenomenology [?, ?] has shown that one can look for quantitative deviations from the necessary pre-conditions of traditional RMT (particularly that the weights are 0-mean and finite variance) to detect when a model layer suffers from some other anomaly in the elements, which we call a “Correlation Trap” (see Section 3.3.1 and [?, ?]). A Correlation Trap may cause, or be caused by, the over-regularization leading Alpha to fall below 2, (see Section 3.3.2). Here, we explore this in greater detail, in light of our SETOL.

Let’s look at the ESDs of the FC1 layer of the MLP3 model, for learning rates  $lr = 16\times$  and  $lr = 32\times$  normal. We will be interested in the general shape of the ESD of  $\frac{1}{N}\mathbf{W}^T\mathbf{W}$ . For the purposes of detecting a Correlation Trap, we will randomize  $\mathbf{W}^T\mathbf{W}$  element-wise, and then observe its largest eigenvalue  $\lambda_{rand}^{max}$ .

Figure 24 shows the ESDs of the original matrix (green) and the element-wise randomized  $\text{rand}(\frac{1}{N}\mathbf{W}^T\mathbf{W})$  (red). Observe in particular  $\lambda_{rand}^{max}$  for each learning rate factor. For  $lr = 16\times$ , Figure 24(a) shows that the ESD of  $\frac{1}{N}\mathbf{W}^T\mathbf{W}$  is HT, whereas the ESD of the randomized matrix is essentially a distorted semi-circle—as expected from the well-known MP result; and that  $\lambda_{rand}^{max}$  lies at the edge of the random MP Bulk ESD. (A similar result is seen for smaller learning rates.) In contrast, for  $lr = 32\times$ , Figure 24(b) shows that while the original ESD is again HT, the ESD of the randomized has one large element,  $\lambda_{rand}^{max}$ , that pulls out from the MP bulk. This is the signature of a Correlation Trap; and it co-occurs with the exact learning rate setting that degraded the train and test accuracies, pushing Alpha below its optimal value of  $\alpha \simeq 2$ . When this happens, both the estimation of Alpha and the formation of a PL tail are potentially disrupted.

Correlation Traps have been observed previously [?, ?], using the HTSR phenomenology. However, SETOL provides an explanation for why this would be expected to occur — non-standard element-wise distributions will tend to interfere with the properties of the spectrum which SETOL analyzes. Our derivation in Section 6 suggests that in order to apply the SETOL effectively one must avoid (or remove) such traps. This too has been observed previously [?, ?].

## 6.5 Overloading and the Hysteresis Effect

Here, we do such-and-such. XXX. HOW PRECISELY TO FRAME.

Obtaining a value of  $\alpha$  outside the  $\alpha \gtrsim 2$  range is indicative of “overloading [?, ?]” that layer. This can be accomplished, e.g., by training only one layer in the MLP3 model. As the two layers

have very different sizes, we see markedly different behaviors, that are nevertheless consistent with theory.

The SETOL theory is based on the idea that NNs undergoing training behave like Statistical Mechanic systems relaxing to an equilibrium. So far, we have tested the theory under conditions that are approaching Ideal. However, for the theory to be useful in practice, we must also examine how it performs in non-Ideal situations. Of particular interest, we would like to examine the theory under conditions where the training dynamics slows down, i.e., when it is in a “glassy or meta-stable state. One way we can do this is to train only one layer, and freeze the rest. Doing so overloads the single trainable weight matrix, as a function of the ratio of examples to trainable parameters  $[\alpha, \beta, \gamma]$ , and we expect this to cause  $\alpha_{FC1}$  or  $\alpha_{FC2}$  to drop well below 2.

### 6.5.1 Baseline: Loading onto both FC1 and FC2

Done

To start, Figure 25 shows  $\alpha_{FC1}$  and  $\alpha_{FC2}$ , binned in units of 0.05, versus train and test error over all epochs of training<sup>42</sup> for each of the different learning rates considered. (Cf. Figures 15 and 16, Section 6.2, which show only the final epoch.) Binning was done so as to facilitate averaging over the 5 starting random seeds; linear fits are shown separately for each learning rate; and error bars represent one standard deviation within each bin. The critical value of  $\alpha = 2$  is shown as a vertical red line in all plots. For each learning rate, train error and  $\alpha$  decrease together during training, which can be seen for both  $\alpha_{FC1}$  (a) and  $\alpha_{FC2}$  (c), reading each line from the top right to bottom left. Test errors, (b, d) show a similar trend, but with wider error bars. Observe that the range of the y-axis is narrower for test error to make detail more visible.

We see in Figure 25 (a,b) that the 32× learning rate causes  $\alpha_{FC1}$  to decrease faster than any other, putting it on course to fall below 2 before train error reaches  $\approx 0$ . (See Figure 14, at the beginning of this Section.) We also see that the slower learning rates cause train error to reach  $\approx 0$  well before  $\alpha_{FC1}$  can reach 2, and this offers an explanation as to why their test error was higher.

### 6.5.2 Overloading FC1

In contrast, when only FC1 is trained, Figure 26 shows that, as expected,  $\alpha_{FC1}$  decreases well below 2. Training error (a) generally trends downward as  $\alpha_{FC1}$  decreases, but no matter the learning rate, the relation is the same, or nearly so, because only one layer is being trained. Consequently, all learning rates were pooled for one linear fit, for visibility.

In Figure 26 we can see demonstration of a crucial claim of SETOL theory: that for  $\alpha_{FC1} > 2$ , (vertical red line,) the test error declines linearly (b) is almost perfectly linear with decline in  $\alpha_{FC1}$ , however, when  $\alpha_{FC1} < 2$ , the curve bends upward. Furthermore, the precision with which the trajectory changes as  $\alpha_{FC1}$  passes the threshold provides ample validation that the estimator of  $[\alpha]$  is indeed accurate. We observe that test error may continue to decrease after  $\alpha_{FC1} < 2$ , however the rate of decrease is significantly less. We can also see that in some sense, the model is “doomed to always have train error reach  $\approx 0$  when  $\alpha_{FC1} \approx 1.7$ , i.e., after  $\alpha = 2$ , because of the number of trainable parameters, and perhaps because of the lack of a modulating influence of FC2 seen in Figure 25.

### 6.5.3 Overloading FC2

The MNIST dataset has 60,000 training examples, which means that an FC1-only model is over-parameterized, but FC2 is substantially *under*-parameterized  $[\alpha]$ . (Table 6) This drastically

<sup>42</sup>Excluding the first four epochs when the matrix is still essentially random.

2372 changes the meaning of an experiment where only FC2 is trained. Figure 25 (Section 6.5.1) shows  
 2373  $\alpha_{FC2}$  vs. train error (c) and test error (d) when *all* layers are trained. There we see a different  
 2374 relationship between  $\alpha_{FC2}$  and train and test error for each learning rate. None of the training  
 2375 runs reached an  $\alpha_{FC2}$  of 2, as the load was split between both layers FC1 and FC2.

2376 Figure 27, however, shows a starkly different relationship. When only FC2 is trained, each  
 2377 random seed produces a different trajectory, meaning that they cannot be plotted as a single curve,  
 2378 even with error bars. Train (a, c) and test (b, d) error rates are shown as a function of  $\alpha_{FC2}$   
 2379 for the two highest learning rate factors,  $16\times$  (a, b) and  $32\times$  (c, d). Lower learning rate factors,  
 2380 (not shown,) showed the same trend, except that they did not progress as far as those shown in  
 2381 Figure 27.

2382 First we see, for all seeds,  $\alpha_{FC2}$  decreases all the way down to 1.5, after which it begins to  
 2383 rebound to the right. As it does, train error continues to decrease. We can also see that test  
 2384 error continues to decrease as well, down to its minimum value of slightly more than 0.03, as  
 2385  $\alpha_{FC2}$  continues to increase for a short time. However, at the exact point where test error reaches  
 2386 a minimum, the PL tail itself begins to fracture, leading to different estimates of  $\alpha_{FC2}$  for each  
 2387 seed. As each of the five starting random seeds are shown separately, we can see that each of them  
 2388 terminates at a different point, some of which are closer to 2 than others.

2389 Such a reversion to a more typical value of  $\alpha_{FC2}$ , prior to the fracturing of the PL tail, is  
 2390 reminiscent of a spin glass system relaxing towards its minimum energy configuration, in a way  
 2391 that retains some memory of the path taken on the way to its current state. We conjecture that  
 2392 if the model were trained for sufficiently many epochs, (perhaps many thousands,) then the tail  
 2393 would re-form, and  $\alpha_{FC2}$  would reform, and revert all the way back to its stable value of 2.



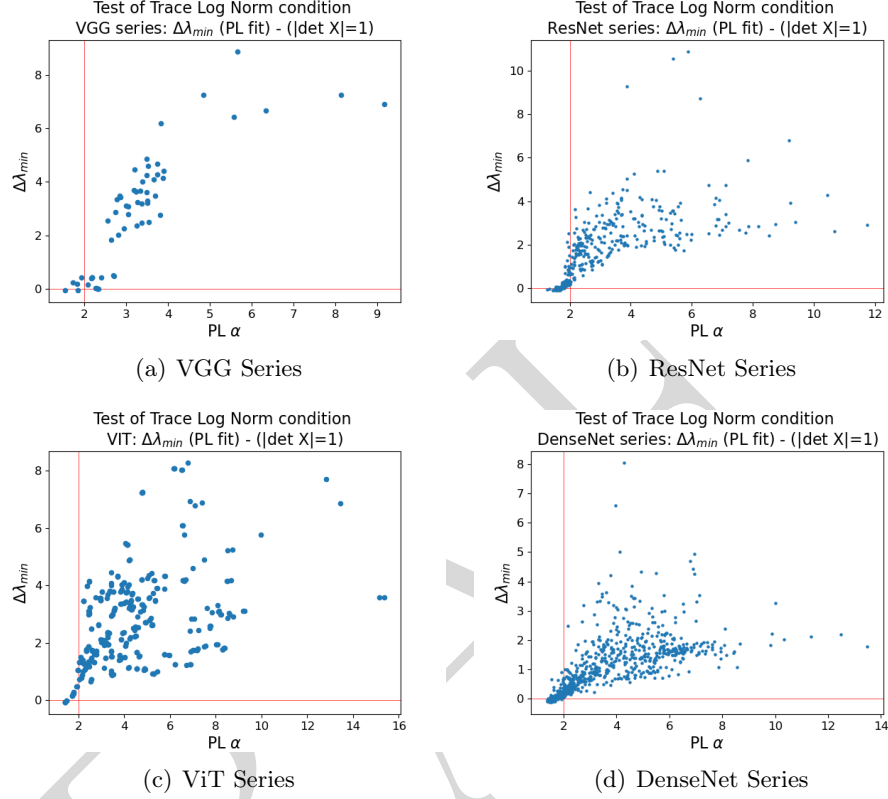


Figure 22: Difference between the two  $\lambda_{min}$  estimates,  $\Delta\lambda_{min}$ , (Eqn. 152), as a function of  $\alpha$ , for linear and convolutional layers in series of VGG [?], ResNet [?], ViT models [?] and DenseNet models [?]. Layer matrices for all models in the series were pooled to create each plot. In (a), (VGG,) we see three clusters of points – those with  $\Delta\lambda_{min}$  close to 0 and  $\alpha$  close to 2, those with  $\Delta\lambda_{min}$  above 2 and  $\alpha > 2.5$ , and those with  $\Delta\lambda_{min}$  above 6 and  $\alpha > 3.5$ . In (b), (ResNet,) we see that in general, as  $\alpha$  shrinks towards 2,  $\Delta\lambda_{min}$  tends towards 0, overshooting slightly. We also see that the difference  $\Delta\lambda_{min}$  is almost always positive, with few exceptions, and even the layers that do not overshoot form a kind of “funnel shape pointing towards the critical point (2,0). In (c), (ViT,) we also see the same general relationship between  $\alpha$  and  $\Delta\lambda_{min}$  across layers of several ViT models. Observe that ViT models do not have convolutional layers, and in spite of this, the overall pattern is similar. In (d), (DenseNet,) we see a similar overall trend as in (b), except that  $\Delta\lambda_{min}$  tends to decrease sooner, but there are also more layers with  $\alpha < 2$  and  $\Delta\lambda_{min}$  above 0.

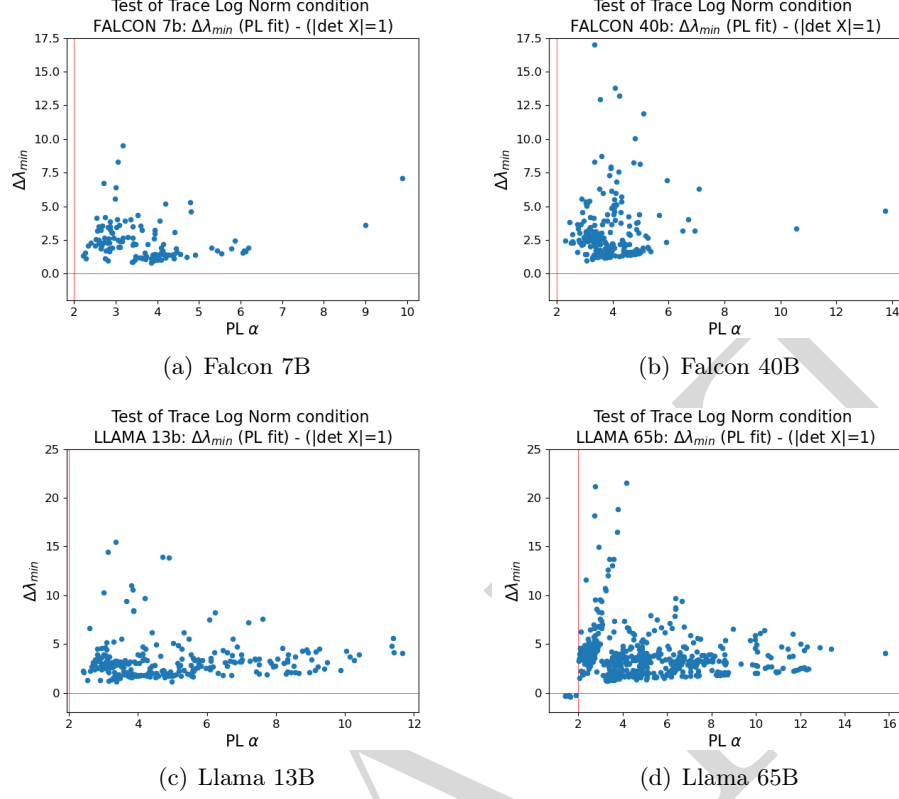


Figure 23: Difference between the two  $\lambda_{min}$  estimates,  $\Delta\lambda_{min} = \lambda_{min}^{PL} - \lambda_{min}^{|\det X|=1}$ , as a function of  $\alpha$ , for all linear layers in the FALCON [?](a-b) and LLAMA [?](c-d) language models for varying numbers of parameters. As in Figure 22, we see that in recent Large Language Models,  $\Delta\lambda_{min}$  remains positive, except where  $\alpha < 2$  (d). Otherwise, a “funnel shape can still be seen leading towards the critical point (2, 0) as in Figures 22(b) and 22(c). Observe that the x- and y-axes are different between sub-figures due to the differences in scale of each model.

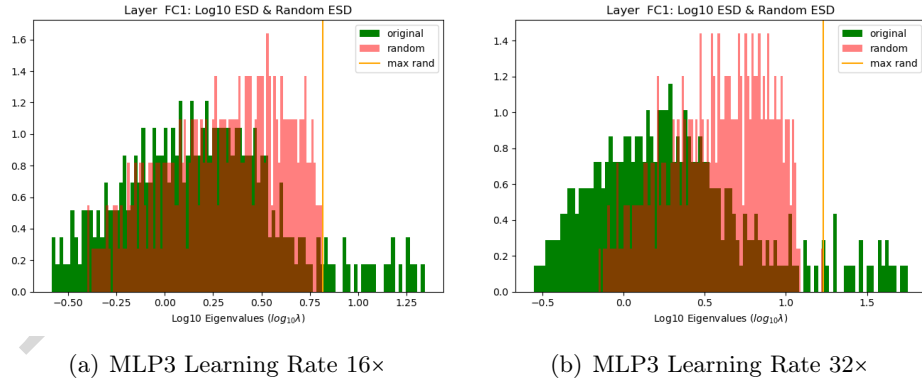


Figure 24: ESD plots for learning rate  $lr = 16\times$  and  $lr = 32\times$  normal, shown on Log-Lin scale, as computed using the **WeightWatcher** tool, for the FC1 weight matrix  $\mathbf{W}$  (green) and an element-wise randomized  $\text{rand}(\mathbf{W})$  (red). This provides an example of inducing a Correlation Trap in the MLP3 model, simply by increasing the learning rate used during model training. See Section 3.3.1.

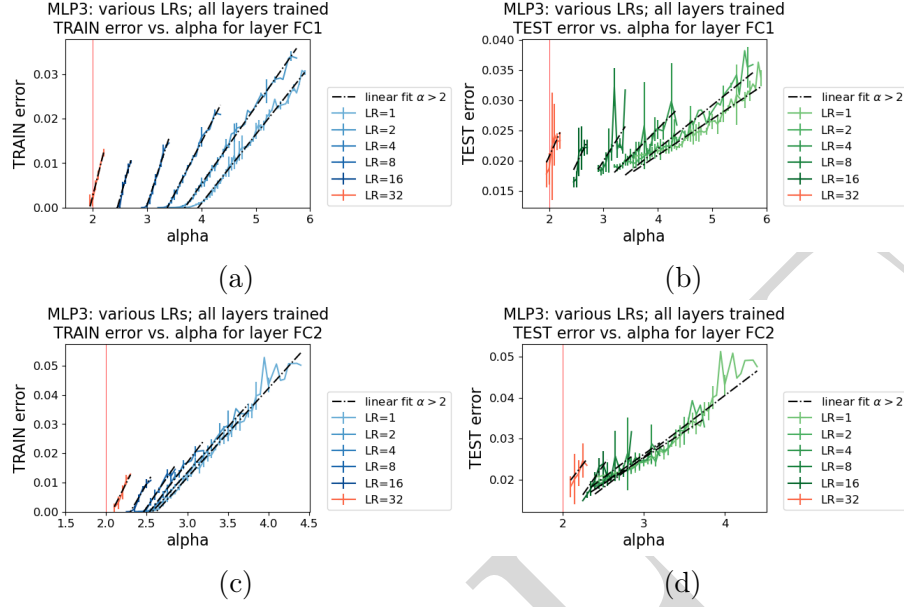


Figure 25: Train (a, c) and test (b, d) accuracy as a function of  $\alpha_{FC1}$  (a, b) and  $\alpha_{FC2}$  (c, d) when **all layers are trained**. Red vertical lines show the critical value of  $\alpha = 2$ , and dashed black lines show linear fits of error, using only points where  $\alpha > 2$  and train error  $> 0.001$ . For FC1 (a, b), we can see that each learning rate produces a different trajectory of train error (a) and test error (b) as a function of  $\alpha_{FC1}$ , showing that even though FC1 dominates overall, (Table 6,) FC2 still plays a modulating role. (Cf. Figure 26 where there is only one trajectory.) In (c, d) we can see that  $\alpha_{FC2}$  never goes below 2. As in (a, b), each learning rate produces a different trajectory, though there is greater overlap for lower learning rates. See discussion in Section 6.5.1.

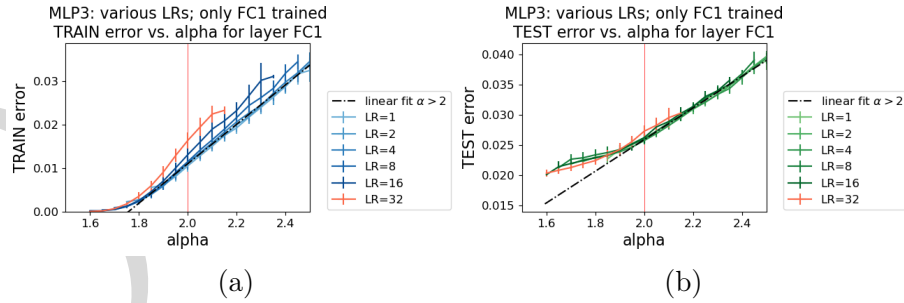
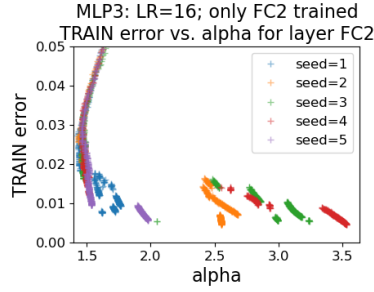
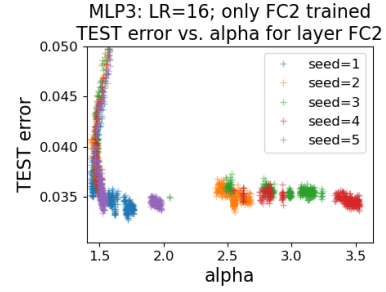


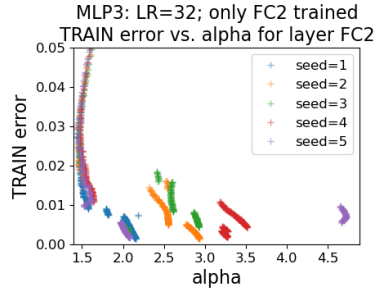
Figure 26: Train and test accuracy as a function of  $\alpha_{FC1}$  when only FC1 is trained (a, b). Red vertical lines show the critical value of  $\alpha_{FC1} = 2$ , and dashed black lines show linear fits of error, using only points where  $\alpha_{FC1} > 2$ . In contrast with Figure 25, when only FC1 is trained, we can see that no matter the learning rate, there is only one trajectory, for both train error (a) and test error (b). Hence, for visibility, only one linear fit, using all learning rates pooled, is shown. Crucially, we can see in (b) that as  $\alpha_{FC1}$  passes below 2, the *test error* trajectory changes, for all learning rates, even as the *train error* trajectory does not, until it reaches  $\sim 0$ . This suggests that even though test accuracy can still decrease when  $\alpha_{FC1} < 2$ , it does so at a decreased rate relative to  $\alpha_{FC1}$ . See discussion in Section 6.5.2.



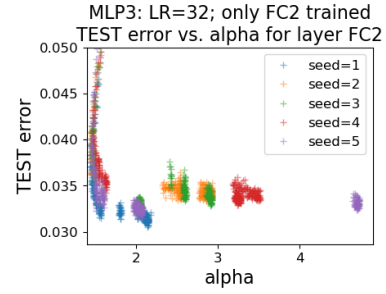
(a)



(b)



(c)



(d)

Figure 27: Train error (a, c) and test error (b, d) as a function of  $\alpha_{FC2}$  when all other layers are frozen, for the two largest learning rates,  $lr = 16\times$ , (a, b) and  $lr = 32\times$  (c, d). Cf. Figure 25, wherein all layers were trained. Due to the markedly different behavior of each random seed, they cannot be plotted as means and error bars, and are instead shown separately. The path taken by all seeds up to  $\alpha_{FC2} = 1.5$  has a slight curvature characteristic of a hysteresis-like behavior. Observe also that the fragmenting into separate paths, due to the breakdown of the PL tail, coincides roughly with each seeds path reaching its minimum test error. The y-axis is scaled differently for train and test error to make variation more visible. See discussion in Section 6.5.3.

## 7 Conclusion and Future Directions

In this work, we have introduced **SETOL**, a *Semi-Empirical Theory of (Deep) Learning* that unifies concepts from Statistical Mechanics (**StatMech**), Heavy-Tailed (HT) Random Matrix Theory (RMT), and quantum-chemistry-inspired approaches to strongly correlated systems [?, ?, ?]. **SETOL** aims to provide a solid theoretical foundation for the Heavy-Tailed Self-Regularization (**HTSR**) phenomenology, including the widely used **Alpha** and **AlphaHat** HT Power Law (PL) Layer Quality metrics, which are implemented in the open-source **WeightWatcher** toolkit. Specifically, **SETOL** reformulates the Neural Network (NN) learning problem for a single layer as a matrix generalization of the classic Student-Teacher (ST) model for Perceptron learning and is analyzed within the Annealed Approximation (AA) at high temperatures (high-T). This reformulation results in a model expressed as an integral over random Student correlation matrices, commonly referred to as the Harish-Chandra–Itzykson–Zuber (HCIZ) integral. To evaluate this integral, we recast the solution using a technique derived from first principles, analogous to a single step of the Exact Wilson Renormalization Group (RG) theory [?]. Leveraging recent results [?, ?], we express the Layer Quality as a sum of integrated matrix-cumulants from RMT (i.e., R-transforms). Finally, we conduct both direct and observational experiments to validate key assumptions of the **SETOL** framework and empirically connect it to the **HTSR** theory.

### Key Contributions and Observations.

- A. *Rigor for the HTSR Phenomenology.* **SETOL** explains *why* power-law (PL) exponents in the layer spectral densities (e.g. **Alpha** and **AlphaHat**) act as robust diagnostics of generalization, even in large, complex architectures without access to training or test data. Our analysis ties these Heavy-Tailed ESDs to a Volume Preserving Transformation associated with an *Effective* Free Energy landscape, and suggesting (in analogy with traditional **StatMech** phases in learning theory) that the **HTSR** condition  $\alpha \approx 2$  marks a phase boundary between optimal generalization and overfitting.
- B. *Matrix-Generalized Student-Teacher (ST) Model.* **SETOL** is formulated as matrix generalization of the classical (vector-based) ST perceptron learning, incorporating  $N \times M$  layer weight matrices,  $\mathbf{w} \rightarrow \mathbf{W}$ . Key to this generalization is isolating the top eigenvalue/eigenvector directions—called the *Effective Correlation Space* (ECS)—before evaluating the resulting partition function (or HCIZ integral). The ECS contains the **HTSR** PL tail, validating that the tail captures the dominant layer generalizing components.
- C. *Trace-Log Condition & **Alpha** = 2.* A remarkable empirical observation, predicted by **SETOL**, is that layers near *ideal* training also satisfy  $\ln(\prod \lambda_i) \approx 0$  in their tail eigenvalues; equivalently,  $\sum \ln \lambda_i \approx 0$ . We call this the **TRACE-LOG** condition (or **DetX** in **WeightWatcher**). Empirically, this condition appears when the **HTSR** **Alpha**  $\approx 2$ . **SETOL** thereby *unifies* two previously separate heuristics for “optimal” or so-called Ideal behavior.
- D. *Empirical Validation of SETOL* To validate the ECS and the **TRACE-LOG** condition, we trained small (3-layer) Multi-Layer Perceptron (MLP3) on MNIST and under varying batch sizes and learning rates. Using this, we verified that when the **HTSR**  $\alpha \approx 2$  the **SETOL** **TRACE-LOG** condition also (usually) holds, and that one can reproduce the training accuracy by retaining only the ECS. This is further confirmed by using the **WeightWatcher** tool to examine the **Alpha** and **DetX** metrics common, open-source CV and NLP models, including modern LLMs (ResNets, DenseNets, ViTs, and LLMs like LLaMA and Falcon).
- E. *Correlation Traps & OverFitting.* We observe that when layer ESDs with  $\alpha < 2$ , they exhibit behavior that can be interpreted as *over-regularization* and/or *overfitting*. For example, we

observe what we call *Correlation Traps*, large rank-one perturbations in the (randomized) layer weight matrix  $\mathbf{W}$  that can be induced by training with excessively small batch sizes (bs=1) and are associated with degraded test accuracy, and which cause the HTSR to drop to  $\alpha < 2$ .

Additionally, we can induce a overfitting by freezing all but one layer and then training, which causes the layer  $\alpha < 2$ . By training in the underparameterized regime, we observe path-dependent, “glassy” behavior.

F. *Connection to Semi-Empirical Methods.* Conceptually, SETOL parallels well-known *Semi-Empirical* methods in quantum chemistry [?, ?, ?], wherein complicated many-body Hamiltonians are approximated by effective theories, but fitted or validated using empirical data. By retaining only the largest spectral modes (ECS) and imposing the TRACE-LOG condition, we can describe crucial low-rank correlations while discarding less relevant interactions, much like in Freed–Martin Effective Hamiltonian theories or Wilson’s Renormalization Group (RG) approach.x

**RG Analogy: A One-Step View.** From a Renormalization Group perspective, restricting to the measure on the partition function to the ECS is akin to performing a *single step* of the Exact Wilson Renormalization Group. In doing this, we are discarding bulk “uninteresting” degrees of freedom in favor of the strongly correlated HT *long-ranged* modes. This leads to an effective model with fewer degrees of freedom but *renormalized* interactions—interactions that are dominated by the largest eigenvalues. This analogy with RG theory suggests that the HTSR phenomenology, where  $\alpha \in [2, 6]$  in the Fat-Tailed Universality Class, is essentially describing a near-critical phase when  $\alpha \approx 2$  and satisfies  $\ln(\prod \lambda_i) \approx 0$  in its ECS. Departing from this point ( $\alpha < 2$  or  $\alpha > 4 - 6$ ) leads to suboptimal results, consistent with the multi-phase pictures in StatMech spin glass theories of learning [?, ?, ?, ?].

**Toward Understanding “Why Deep Learning Works.”** A key question in deep learning theory is why large neural networks achieve strong generalization despite operating in highly non-convex optimization landscapes. From the perspective of RG exact theory, this phenomenon can be partially understood through the concentration of generalization-relevant components. Specifically, models trained in regimes exhibiting *Lévy-like* or *Power-Law (PL)* couplings lead to the emergence of effective low-dimensional descriptions, where irrelevant modes are suppressed. The **WeightWatcher Alpha** metric quantitatively captures this concentration by characterizing how the optimization landscape stabilizes near criticality, with  $\alpha \approx 2$  signaling optimal generalization in a near-critical regime.

**Relation to Levy Spin Glasses and Heavy-Tailed Random Matrix Models.** One long-standing puzzle is why large NNs avoid the worst of highly non-convex optimization, despite nominal exponential degeneracies. SETOL offers a partial explanation: if the trained model has *Lévy-like* or *PL* couplings, then typical spin-glass degeneracies can be lifted, leaving the layer in a finite number of near-critical minima. In analogy with older work on *Levy Spin Glasses*[?], it is proposed **WeightWatcher Alpha** metric effectively measures how *rugged yet stable* the effective energy landscape is for each layer, with  $\alpha \approx 2$  signifying a sweet spot of Ideal generalization.

## 7.1 Future Directions

**1. Multi-Layer RG and Layer Interactions.** While SETOL is formulated per-layer, modern DNNs stack many layers, each potentially with different  $\alpha$ . An improved approach would treat



layer-layer interactions, potentially with a multi-step RG or other approach, that could track how these exponents evolve or couple across depth. It is possible that certain layers (e.g. final fully-connected heads in LLMs) exhibit  $\alpha$  far from 2, while others converge near 2—raising questions about how best to address or combine them.

**2. Explicit R-transform Expansions for Heavy-Tailed ESDs.** In practice, large pre-trained models often show truncated power-law (TPL) ESDs, with  $\alpha \geq 2$ , and with sharp tail cutoff (e.g.  $\lambda_{\max}$ ) [?, ?]. A possible extension of SETOL could treat these cases by deriving the explicit form of the R-transform for PL/TPL ESDs with  $\alpha \in [2, 3, 4]$ , and applying it even in cases with the TRACE-LOG condition does not strictly hold.

**3. Practical Diagnostics and Fine-Tuning.** The open-source WeightWatcher tool has already seen success diagnosing layer quality. Integrating SETOL’s TRACE-LOG condition may refine this further, helping users identify correlation traps or “under-exploited” layers. There is also strong potential for using Alpha or TRACE-LOG-based signals during training or fine-tuning: e.g. automatically adapting learning rates to push each layer closer to  $\alpha = 2$ , for fine-tuning models with significantly less memory[?], for compressing large LLMs[?], and other practical applications.

**4. Correlation Traps and Meta-Stable States.** Although our experiments show how small batch sizes or large learning rates can induce correlation traps, a quantitative theory of *where and why* traps occur remains open. Clarifying these states could enable *trap-avoidance* strategies, e.g. partial re-initialization or specialized regularizers that favor lower-rank updates in the ECS. In large-language-model (LLM) contexts, correlation traps might manifest as *hallucinations* or *mode collapse*, motivating deeper analysis.

**4. Analyzing the Layer Null Space.** One critical but often overlooked factor is the potential null space within model layers, which can emerge during overfitting. This null space represents parameter directions that fail to contribute meaningfully to generalization but instead encode redundant or overly specific patterns tied to the training data which might be ignored or forgotten. Future work should examine if and when NN layers have components in their null space, which contribute significantly to the performance of the model.

**5. Layer-Layer Cross-Terms.** SETOL is a single layer theory, however, as noted in Section 5.1, it would be desirable to extend the theory to including layer-layer cross terms. While we don’t have an exact expression for this, we can propose a phenomenological guess that the leading order term would be the integrated R-transform, defined for the overlap between nearest-neighbor weight matrices (i.e.,  $\mathbf{W}_1, \mathbf{W}_2$ ) that can be aligned along a common axis. This term would take the form  $\mathbb{G}_{\mathbf{R}_{1,2}}(\lambda_{1,2})$  where  $\mathbf{R}_{1,2} \sim \mathbf{W}_1^\top \mathbf{W}_2$  and  $\lambda_{1,2}$  is an eigenvalue of  $\mathbf{R}_{1,2}$ . Note that the open-source WeightWatcher tool can identify and compute the intra-layer interactions.[?]

**Concluding Remarks.** SETOL as a Semi-Empirical theory merges first-principles methods from StatMech and RMT with empirical insights from HTSR and the open-source WeightWatcher tool. It clarifies *how* Heavy-Tailed layer weight matrices can emerge from training on realistic data, and *why* their spectral exponents so reliably predict generalization quality without peeking at training/test sets. In so doing, SETOL not only offers new insights into the “*why does it work*” question of deep learning but also suggests a roadmap for improving DNN models by focusing attention on that near-critical subspace of their largest eigenvalues. We are optimistic that future developments along these lines—extending the single-step Renormalization Group analogy, refining

TPL expansions, and systematically diagnosing correlation traps—will yield more robust, data-free metrics for training, fine-tuning, and compressing next-generation neural networks.

**Acknowledgements.** We would like to acknowledge Matt Lee of Triaxiom Capital and Carl Page of the Anthropocene Institute. We also thank Mirco Milletari for helpful conversations. MWM would like to acknowledge DARPA, NSF, and ONR as well as the UC Berkeley BDD project and a gift from Intel for providing partial support of this work. Our conclusions do not necessarily reflect the position or the policy of our sponsors, and no official endorsement should be inferred. XXX. WHO ELSE TO THANK.

## References

- [1] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Židek, Anna Potapenko, et al. Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589, 2021.
- [2] Nobel Prize Organization. The nobel prize in physics 2024, 2024.
- [3] Nobel Prize Organization. The nobel prize in chemistry 2024, 2024.
- [4] A. Engel. Complexity of learning in artificial neural networks. *Theoretical Computer Science*, 265(1–2):285–306, 2001.
- [5] A. Engel and C. P. L. Van den Broeck. *Statistical mechanics of learning*. Cambridge University Press, New York, NY, USA, 2001.
- [6] E Gardner. The space of interactions in neural network models. *Journal of Physics A: Mathematical and General*, 21(1):257, jan 1988.
- [7] H. Sompolinsky, N. Tishby, and H. S. Seung. Learning from examples in large neural networks. *Phys. Rev. Lett.*, 65:1683–1686, Sep 1990.
- [8] H. S. Seung, H. Sompolinsky, and N. Tishby. Statistical mechanics of learning from examples. *Physical Review A*, 45(8):6056–6091, 1992.
- [9] E. Levin, N. Tishby, and S. A. Solla. A statistical approach to learning and generalization in layered neural networks. *Proceedings of the IEEE*, 78(10):1568–1574, 1990.
- [10] Erin Grant, Sandra Nestler, Berfin Şimşek, and Sara Solla. Statistical physics, Bayesian inference and neural information processing. *arXiv e-prints*, page arXiv:2309.17006, September 2023.
- [11] V.N. Vapnik. *Statistical Learning Theory*. John Wiley & Sons, New York, 1998.
- [12] J. J. Hopfield. Neural networks and physical systems with emergent collective computational abilities. *Proc. Natl. Acad. Sci. USA*, 79(8):2554–2558, 1982.
- [13] D. H. Ackley, G. E. Hinton, and T. J. Sejnowski. A learning algorithm for Boltzmann machines. *Cognitive Science*, 9(1):147–169, 1985.
- [14] G. E. Hinton and T. J. Sejnowski. Learning and relearning in Boltzmann machines. In D. E. Rumelhart, J. L. McClelland, and CORPORATE PDP Research Group, editors, *Parallel distributed processing: explorations in the microstructure of cognition, vol. 1*, pages 282–317. MIT Press, 1986.
- [15] W. A. Little. The existence of persistent states in the brain. *Math. Biosci.*, 19:101–120, 1974.
- [16] M. Belkin, D. Hsu, S. Ma, and S. Mandal. Reconciling modern machine-learning practice and the classical bias–variance trade-off. *Proc. Natl. Acad. Sci. USA*, 116:15849–15854, 2019.
- [17] Marco Loog, Tom Viering, Alexander Mey, and David MJ Tax. A brief prehistory of double descent. *Proceedings of the National Academy of Sciences*, 117(20):10625, 2020.
- [18] M. Oppen. Learning to generalize. In D. Baltimore, editor, *Frontiers of Life: Intelligent Systems*, pages 763–775. Academic Press, Cambridge, 2001.

- [19] David A. Roberts, Sho Yaida, and Boris Hanin. *The Principles of Deep Learning Theory: An Effective Theory Approach to Understanding Neural Networks*. Cambridge University Press, 2022.
- [20] V. Vapnik, E. Levin, and Y. Le Cun. Measuring the VC-dimension of a learning machine. *Neural Computation*, 6(5):851–876, 1994.
- [21] T. L. H. Watkin, A. Rau, and M. Biehl. The statistical mechanics of learning a rule. *Rev. Mod. Phys.*, 65(2):499–556, 1993.
- [22] D. Haussler, M. Kearns, H. S. Seung, and N. Tishby. Rigorous learning curve bounds from statistical mechanics. *Machine Learning*, 25(2):195–236, 1996.
- [23] G. K. Dziugaite and D. M. Roy. Computing nonvacuous generalization bounds for deep (stochastic) neural networks with many more parameters than training data. Technical Report Preprint: arXiv:1703.11008, 2017.
- [24] C. H. Martin and M. W. Mahoney. Post-mortem on a deep learning contest: a Simpson’s paradox and the complementary roles of scale metrics versus shape metrics. Technical Report Preprint: arXiv:2106.00734, 2021.
- [25] C. H. Martin and M. W. Mahoney. Implicit self-regularization in deep neural networks: Evidence from random matrix theory and implications for learning. *Journal of Machine Learning Research*, 22(165):1–73, 2021.
- [26] C. H. Martin, T. S. Peng, and M. W. Mahoney. Predicting trends in the quality of state-of-the-art neural networks without access to training or testing data. *Nature Communications*, 12(4122):1–13, 2021.
- [27] Yaoqing Yang, Ryan Theisen, Liam Hodgkinson, Joseph E. Gonzalez, Kannan Ramchandran, Charles H. Martin, and Michael W. Mahoney. Evaluating natural language processing models with generalization metrics that do not need access to any training or testing data. Technical Report Preprint: arXiv:2202.02842, 2022.
- [28] Y. Yang, R. Theisen, L. Hodgkinson, J. E. Gonzalez, K. Ramchandran, C. H. Martin, and M. W. Mahoney. Test accuracy vs. generalization gap: Model selection in NLP without accessing training or testing data. In *Proceedings of the 29th Annual ACM SIGKDD Conference*, pages 3011–3021, 2023.
- [29] Yefan Zhou, TIANYU PANG, Keqin Liu, Charles Martin, Michael W Mahoney, and Yaoqing Yang. Temperature balancing, layer-wise weight analysis, and neural network training. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 63542–63572. Curran Associates, Inc., 2023.
- [30] Haiquan Lu, Yefan Zhou, Shiwei Liu, Zhangyang Wang, Michael W. Mahoney, and Yaoqing Yang. Alphapruning: Using heavy-tailed self regularization theory for improved layer-wise pruning of large language models. In *Advances in Neural Information Processing Systems*, volume 37, pages 0000–0000, 2024.
- [31] John Negele. Hans Bethe and the theory of nuclear matter. *Physics Today*, 58(10):58, 2005.
- [32] D. Ivanenko. The proton-neutron hypothesis of atomic nuclei. *Nature*, 129:798, 1932.
- [33] Maria Goeppert-Mayer. On closed shells in nuclei. ii. *Physical Review*, 75(10):1969–1970, 1949.
- [34] J. Hans D. Jensen, Otto Haxel, and Hans Suess. On the “magic numbers” in nuclear structure. *Physical Review*, 75:1766, 1949.
- [35] Eugene Wigner. Characteristic vectors of bordered matrices with infinite dimensions. *Annals of Mathematics*, 62(3):548–564, 1955.
- [36] V. A. Marchenko and L. A. Pastur. Distribution of eigenvalues for some sets of random matrices. *Mathematics of the USSR-Sbornik*, 1(4):457–483, 1967.
- [37] T. Guhr, A. Müller-Groeling, and H. A. Weidenmüller. Random matrix theories in quantum physics: Common concepts. *Physics Reports*, 299:190, 1998.

- [38] A. Zee. Law of addition in random matrix theory. *Nuclear Physics B*, 474(3):726–744, September 1996.
- [39] Melih K. Sener and Klaus Schulten. General random matrix approach to account for the effect of static disorder on the spectral properties of light harvesting systems. *Physical Review E*, 65(3):031916, 2002. Received 6 June 2001; revised manuscript received 29 August 2001; published 6 March 2002.
- [40] S. Galluccio, J.-P. Bouchaud, and M. Potters. Rational decisions, random matrices and spin glasses. *Physica A*, 259:449–456, 1998.
- [41] R. Cherrier, D. S. Dean, and A. Lefèvre. Role of the interaction matrix in mean-field spin glass models. *Physical Review E*, 67(4), April 2003.
- [42] Rudolph Pariser and Robert G. Parr. A semi-empirical theory of the electronic spectra and electronic structure of complex unsaturated molecules. i. *The Journal of Chemical Physics*, 21(3):466–471, 1953.
- [43] J. Hubbard. Electron correlations in narrow energy bands. *Proceedings of the Royal Society of London. Series A, Mathematical and Physical Sciences*, 276(1365):238–257, 1963.
- [44] Michael J. S. Dewar and Walter Thiel. Ground states of molecules. 38. the mindo/3 method. approximations and parameters. *Journal of the American Chemical Society*, 97(16):4899–4907, 1975.
- [45] John Ridley and Michael C. Zerner. Intermediate neglect of differential overlap spectroscopy: a reexamination using a modified neglect of differential overlap approach. *Theoretica Chimica Acta*, 32:111–134, 1973.
- [46] James J. P. Stewart. Mopac: A semiempirical molecular orbital program. *Journal of Computer-Aided Molecular Design*, 4:1–103, 1990.
- [47] Arieh Warshel and Michael Levitt. Theoretical studies of enzymic reactions: dielectric, electrostatic and steric stabilization of the carbonium ion in the reaction of lysozyme. *Journal of Molecular Biology*, 103(2):227–249, 1976.
- [48] J. Hubbard. Calculation of partition functions. *Physical Review Letters*, 3(2):77–78, 1959.
- [49] K. F. Freed. Theoretical basis for semiempirical theories. In G.A. Segal, editor, *Semiempirical Methods of Electronic Structure Calculation*, volume 7 of *Modern Theoretical Chemistry*. Springer, 1977.
- [50] Karl F. Freed. Is there a bridge between ab initio and semiempirical theories of valence? *Accounts of Chemical Research*, 16:137–144, Mar 1983.
- [51] Charles H. Martin and Karl F. Freed. Ab initio computation of semiempirical  $\pi$ -electron methods. v. geometry dependence of  $h\nu$   $\pi$ -electron effective integrals. *The Journal of Chemical Physics*, 105(4):1437–1450, 1996.
- [52] Charles H. Martin. Highly accurate ab initio  $\pi$ -electron hamiltonians for small protonated schiff bases. *The Journal of Physical Chemistry*, 100:14310–14315, 1996.
- [53] Charles H Martin. Redesigning semiempirical-like  $\pi$ -electron theory with second order effective valence shell hamiltonian (hv) theory: application to large protonated schiff bases. *Chemical Physics Letters*, 257(3-4):229–237, 1996.
- [54] Charles H. Martin and Robert R. Birge. Reparametrizing mndo for excited-state calculations by using ab initio effective hamiltonian theory: Application to the 2,4-pentadien-1-iminium cation. *The Journal of Physical Chemistry A*, 102(5):852–860, 1998.
- [55] Nobel Prize Committee. The nobel prize in physics 1982: Kenneth g. wilson. <https://www.nobelprize.org/prizes/physics/1982/wilson/>, 1982. Accessed: 2024-12-09.
- [56] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training. *OpenAI*, 2018.
- [57] J. M. Jumper, K. F. Freed, and T. R. Sosnick. Maximum-likelihood, self-consistent side chain free energies with applications to protein molecular dynamics. Technical Report Preprint: arXiv:1610.07277, 2016.

- [58] D. A. Roberts, S. Yaida, and B. Hanin. *The Principles of Deep Learning Theory*. Cambridge University Press, 2021.
- [59] J. Lee, Y. Bahri, R. Novak, S. S. Schoenholz, J. Pennington, and J. Sohl-Dickstein. Deep neural networks as Gaussian processes. Technical Report Preprint: arXiv:1711.00165, 2017.
- [60] G. Yang. Tensor programs III: Neural matrix laws. Technical Report Preprint: arXiv:2009.10685, 2021.
- [61] C. H. Martin and M. W. Mahoney. Traditional and heavy-tailed self regularization in neural network models. In *Proceedings of the 36th International Conference on Machine Learning*, pages 4284–4293, 2019.
- [62] C. H. Martin and M. W. Mahoney. Heavy-tailed Universality predicts trends in test accuracies for very large pre-trained deep neural networks. In *Proceedings of the 20th SIAM International Conference on Data Mining*, 2020.
- [63] B. Derrida. Random-energy model: An exactly solvable model of disordered systems. *Physical Review B*, 24:2613–2626, Sep 1981.
- [64] Joseph D. Bryngelson and Peter G. Wolynes. Spin glasses and the statistical mechanics of protein folding. *Proceedings of the National Academy of Sciences of the United States of America*, 84:7524–7528, Nov 1987.
- [65] Charles H. Martin. Weightwatcher, 2023. Accessed: 2024-06-03.
- [66] A. Clauset, C. R. Shalizi, and M. E. J. Newman. Power-law distributions in empirical data. *SIAM Review*, 51(4):661–703, 2009.
- [67] J. Alstott, E. Bullmore, and D. Plenz. powerlaw: A python package for analysis of heavy-tailed distributions. *PLoS ONE*, 9(1):e85777, 2014.
- [68] Matthias Thamm, Max Staats, and Bernd Rosenow. Random matrix analysis of deep neural network weight matrices. *Physical Review E*, 106(5):054124, 2022.
- [69] J.-P. Bouchaud and M. Potters. *Theory of Financial Risk and Derivative Pricing: From Statistical Physics to Risk Management*. Cambridge University Press, 2003.
- [70] A. Edelman and Y. Wang. Random matrix theory and its innovative applications. In R. Melnik and I. Kotsireas, editors, *Advances in Applied Mathematics, Modeling, and Computational Science*. Springer, 2013.
- [71] Marc Potters and Jean-Philippe Bouchaud. *A First Course in Random Matrix Theory: for Physicists, Engineers and Data Scientists*. Cambridge University Press, 2020.
- [72] Francis R Bach and Michael I Jordan. Learning spectral clustering, with application to speech separation. *The Journal of Machine Learning Research*, 7:1963–2001, 2006.
- [73] Florent Benaych-Georges and Raj Rao Nadakuditi. The eigenvalues and eigenvectors of finite, low rank perturbations of large random matrices. *Advances in Mathematics*, 227(1):494–521, 2011.
- [74] Per Bak. *How nature works: the science of self-organized criticality*. Oxford University Press, Oxford, UK, 1997.
- [75] D. Sornette. *Critical phenomena in natural sciences: chaos, fractals, selforganization and disorder: concepts and tools*. Springer-Verlag, Berlin, 2006.
- [76] L. Laloux, P. Cizeau, J.-P. Bouchaud, and M. Potters. Noise dressing of financial correlation matrices. *Phys. Rev. Lett.*, 83(7):1467–1470, 1999.
- [77] L. Laloux, P. Cizeau, M. Potters, and J.-P. Bouchaud. Random matrix theory and financial correlations. *Mathematical Models and Methods in Applied Sciences*, pages 109–11, 2005.
- [78] Y. Yang, L. Hodgkinson, R. Theisen, J. Zou, J. E. Gonzalez, K. Ramchandran, and M. W. Mahoney. Taxonomizing local versus global structure in neural network loss landscapes. Technical Report Preprint: arXiv:2107.11228, 2021.

- [79] C. H. Martin and M. W. Mahoney. Heavy-tailed Universality predicts trends in test accuracies for very large pre-trained deep neural networks. Technical Report Preprint: arXiv:1901.08278, 2019.
- [80] H. Sompolinsky, N. Tishby, and H. S. Seung. Learning from examples in large neural networks. *Phys. Rev. Lett.*, 65:1683–1686, 1990.
- [81] T. Tanaka. On dualistic structure involving shannon transform and integrated r-transform. pages 1651–1654, 2007.
- [82] T. Tanaka. Asymptotics of Harish-Chandra-Itzykson-Zuber integrals and free probability theory. *J. Phys.: Conf. Ser.*, 95(1):012002, 2008.
- [83] M. Gurbuzbalaban, U. Simsekli, and L. Zhu. The heavy-tail phenomenon in SGD. Technical Report Preprint: arXiv:2006.04740, 2020.
- [84] Chaim Baskin, Evgenii Zheltonozhkii, Tal Rozen, Natan Liss, Yoav Chai, Eli Schwartz, Raja Giryes, Alexander M Bronstein, and Avi Mendelson. Nice: Noise injection and clamping estimation for neural network quantization. *Mathematics*, 9(17):2144, 2021.
- [85] Jungwook Choi, Zhuo Wang, Swagath Venkataramani, Pierce I-Jen Chuang, Vijayalakshmi Srinivasan, and Kailash Gopalakrishnan. Pact: Parameterized clipping activation for quantized neural networks. *arXiv preprint arXiv:1805.06085*, 2018.
- [86] A. Engel, C. Van den Broeck, and C. Broeck. *Statistical Mechanics of Learning*. Statistical Mechanics of Learning. Cambridge University Press, 2001.
- [87] Hanoch Gutfreund, Haim Sompolinsky, and Daniel Stein. Spin-glass models of neural networks. *Physical Review A*, 32(2):1007–1018, 1985.
- [88] C. H. Martin and M. W. Mahoney. Rethinking generalization requires revisiting old ideas: statistical mechanics approaches and complex learning behavior. Technical Report Preprint: arXiv:1710.09553, 2017.
- [89] B. H. Brandow. Foundations of the nuclear shell model. *Physics Letters*, 4(8):294–296, 1963.
- [90] C. H. Martin and M. W. Mahoney. Rethinking generalization requires revisiting old ideas: statistical mechanics approaches and complex learning behavior. Technical Report Preprint: arXiv:1710.09553v1, 2017.
- [91] G. Parisi and M. Potters. Mean-field equations for spin models with orthogonal interaction matrices. *Journal of Physics A: Mathematical and General*, 28(18):5267–5286, 1995.
- [92] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals. Understanding deep learning requires rethinking generalization. Technical Report Preprint: arXiv:1611.03530, 2016.
- [93] Bradley Efron and Robert J Tibshirani. *An Introduction to the Bootstrap*. Chapman & Hall/CRC, 1993.
- [94] F. Rosenblatt. *Principles of Neurodynamics*. Spartan, New York, NY, USA, 1962.
- [95] Andrew M Saxe, James L McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. *arXiv preprint arXiv:1312.6120*, 2013.
- [96] Joël Bun. *Application of Random Matrix Theory to High Dimensional Statistics*. Phd thesis, Université Paris Saclay (COMUE), 2016. NNT: 2016SACLS245, tel-01400544.
- [97] Z. Burda, J. Jurkiewicz, M. A. Nowak, G. Papp, and I. Zahed. Lévy matrices and financial covariances. Technical Report Preprint: arXiv:cond-mat/0103108, 2001.
- [98] Z. Burda, J. Jurkiewicz, M. A. Nowak, G. Papp, and I. Zahed. Random Lévy matrices revisited. Technical Report Preprint: arXiv:cond-mat/0602087, 2006.
- [99] Z. Burda and J. Jurkiewicz. Heavy-tailed random matrices. Technical Report Preprint: arXiv:0909.5228, 2009.



- [100] Jipeng Li, Xueqiong Yuan, and Ercan Engin Kuruoglu. Exploring weight distributions and dependence in neural networks with  $\alpha$ -stable distributions. *IEEE Transactions on Artificial Intelligence*, 5(11):5519–5535, November 2024.
- [101] LeCun Yann. The mnist database of handwritten digits. *R*, 1998.
- [102] H. Nishimori. *Statistical Physics of Spin Glasses and Information Processing: An Introduction*. Oxford University Press, Oxford, 2001.
- [103] S. L. Smith, P.-J. Kindermans, C. Ying, and Q. V. Le. Don’t decay the learning rate, increase the batch size. Technical Report Preprint: arXiv:1711.00489, 2017.
- [104] M. Welling and Y. W. Teh. Bayesian learning via stochastic gradient Langevin dynamics. In *Proceedings of the 28th International Conference on Machine Learning*, pages 681–688, 2011.
- [105] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. Technical Report Preprint: arXiv:1512.03385, 2015.
- [106] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. Technical Report Preprint: arXiv:1409.1556, 2014.
- [107] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [108] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708, 2017.
- [109] Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Merouane Debbah, Etienne Goffinet, Daniel Heslow, Julien Launay, Quentin Malartic, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. Falcon-40B: an open large language model with state-of-the-art performance. 2023.
- [110] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutu Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [111] M. Dereziński, F. Liang, and M. W. Mahoney. Exact expressions for double descent and implicit regularization via surrogate random design. Technical Report Preprint: arXiv:1912.04533, 2019.
- [112] Peijun Qing, Chongyang Gao, Yefan Zhou, Xingjian Diao, Yaoqing Yang, and Soroush Vosoughi. Alphasora: Assigning lora experts based on layer training quality, 2024.
- [113] François Chollet. keras. <https://github.com/fchollet/keras>, 2015.
- [114] X. Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the 13th International Workshop on Artificial Intelligence and Statistics*, pages 249–256, 2010.
- [115] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- [116] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [117] A. Engel and C. Van den Broeck. *Statistical Mechanics of Learning*. Cambridge University Press, Cambridge, UK, 2001.
- [118] R. Cherrier, D. S. Dean, and A. Lefèvre. Role of the interaction matrix in mean-field spin glass models. *Physical Review E*, 67:046112, 2003.

## A Appendix

### A.1 Data Vectors, Weight Matrices, and Other Symbols

See Table 7 for a summary of various vectors and matrices, including their dimensions; see Table 8 for summary of various various symbols used throughout the text; and see Table 9 for summary of types of “Energies” used throughout the text.

Number of NN Layers	index $L$	$N_L$
Number of Data Examples	index $\mu$	$N = N_D$
Number of (input) Features	index $i, j$	$N_f$
Actual Data (Matrix)	$D$	$N_f \times N_D$
Model Data (Matrix)	$\mathcal{D}$	$N_f \times N$
Teacher Perceptron Weight Vector	$\mathbf{t}$	$m$
Student Perceptron Weight Vector	$\mathbf{s}$	$m$
Actual Input Data Vector	$\mathbf{x}_\mu$	$N_f \times 1$
Gaussian model of Input Data Vector	$\boldsymbol{\xi}_\mu$	$N_f \times 1$
Actual Input Data Label	$y_\mu$	$+1 -1$
Model Input Data Label	$y_\mu$	$+1 -1$
General Weight Matrix	$\mathbf{W}$	$N \times M$
General Correlation Matrix	$\mathbf{X} = \frac{1}{N} \mathbf{W}^\top \mathbf{W}$	$M \times M$
Input Layer Weight Matrix	$\mathbf{W}_1$	$N \times M$
Hidden Layer Weight Matrix	$\mathbf{W}_2$	$N \times M$
Output Layer Weight Matrix	$\mathbf{W}_3$	$M \times 2$
Teacher Weight Matrix	$\mathbf{T}$	$N \times M$
Student Weight Matrix	$\mathbf{S}$	$N \times M$
Student-Teacher Overlap Matrix	$\mathbf{R} = \frac{1}{N} \mathbf{S}^\top \mathbf{T}$	$M \times M$
Student Correlation Matrix 1	$\mathbf{A}_1 = \frac{1}{N} \mathbf{S}^\top \mathbf{S}$	$M \times M$
Student Correlation Matrix 2	$\mathbf{A}_2 = \frac{1}{N} \mathbf{S} \mathbf{S}^\top$	$N \times N$
ECS Student Correlation Matrix 1	$\tilde{\mathbf{A}}_1$	$\tilde{M} \times \tilde{M}$
ECS Student Correlation Matrix 2	$\tilde{\mathbf{A}}_2$	$N \times N$
ECS Teacher Correlation Matrix	$\tilde{\mathbf{X}}$	$\tilde{M} \times \tilde{M}$

Table 7: Summary of of various vectors and matrices, including their dimensions.

Perceptron Student-Teacher (ST) Overlap	$R = \frac{1}{N} \mathbf{s}^\top \mathbf{t} = \frac{1}{N} \sum_i s_i t_i$
Student-Teacher (ST) Overlap Operator	$\mathbf{R} = \frac{1}{N} \mathbf{S}^\top \mathbf{T}$
Matrix Generalized ST Overlap	$\frac{1}{N^2} \text{Tr} [\mathbf{R}^\top \mathbf{R}]$
Student-Teacher Self-Overlap	$\eta(\boldsymbol{\xi}) = \mathbf{y}_T^\top \mathbf{y}_S$
$\ell_2$ -Energy Difference or $\ell_2$ -Error	$\Delta E_{\ell_2}(\mathbf{s}, \mathbf{t}, \boldsymbol{\xi})$
$\ell_2$ -Energy Difference Operator Form	$\Delta \mathbf{E}_{\ell_2}(\mathbf{s}, \mathbf{t}, \boldsymbol{\xi}^n) := \sum_{\boldsymbol{\xi}} \Delta E_{\ell_2}(\mathbf{s}, \mathbf{t}, \boldsymbol{\xi})$
$\ell_2$ -Energy Difference Matrix Operator Form	$\Delta \mathbf{E}_{\ell_2}(\mathbf{S}, \mathbf{T}) := \mathbf{I}_M - \frac{1}{N} \mathbf{S}^\top \mathbf{T}$
Annealed Error Potential	$\epsilon(R) = \epsilon(S, T) = \langle \Delta \mathbf{E}_{\ell_2}(\mathbf{s}, \mathbf{t}, \boldsymbol{\xi}^n) \rangle_{\boldsymbol{\xi}^N}$
Linear Perceptron $\epsilon(R)$ at high-T, large- $N$	$\epsilon(R) = 1 - R$
Annealed Approximation (AA)	$\langle \ln Z \rangle_{\boldsymbol{\xi}^N} = \ln \langle Z \rangle_{\boldsymbol{\xi}^N}$
Annealed Hamiltonian	$H^{an}(\mathbf{w})$
Annealed Hamiltonian at high-T	$H_{hT}^{an}(\mathbf{w}) = \epsilon(\mathbf{w})$
Average Student-Teacher Generalization Error	$\bar{\mathcal{E}}_{gen}^{ST} = \langle \epsilon(R) \rangle_{\mathbf{s}}^\beta$
Average Student-Teacher Generalization Accuracy	$1 - \bar{\mathcal{E}}_{gen}^{ST} = \langle \eta(R) \rangle_{\mathbf{s}}^\beta$
Matrix Layer Quality-Squared	$\bar{Q}^2 = \langle \mathbf{R}^\top \mathbf{R} \rangle_{\mathbf{S}}^\beta$
Equivalent Notation for Averages	$\langle \cdots \rangle_A = \int \cdots d\mu(\mathbf{A}) = \mathbb{E}_{\mathbf{A}}[\cdots]$
Projection Operator onto ECS	$\mathbf{P}^{ecs} := \sum  \tilde{\lambda}_i\rangle \langle \tilde{\lambda}_i , i = 1 \cdots \tilde{M}$
Average over ECS Student Correlation Matrices	$\langle \cdots \rangle_{\tilde{\mathbf{A}}}^\beta = \int \cdots d\mu(\tilde{\mathbf{A}}) = \mathbb{E}_{\tilde{\mathbf{A}}}[\cdots]$
TRACE-LOG Determinant Relation	$\text{Tr} [\ln \mathbf{A}] = \ln \det \mathbf{A}$
TRACE-LOG (Volume Preserving) Condition	$\text{Tr} [\ln \mathbf{A}] = 0$ or $\det \mathbf{A} = 1$
Effective Correlation Measure Transform	$d\mu(\mathbf{S}) \rightarrow d\mu(\tilde{\mathbf{A}})$
HCIZ Integral (Tanaka's Notation)	$\mathbb{E}_{\tilde{\mathbf{A}}}[\exp(\beta \text{Tr} [\mathbf{T}^\top \mathbf{A} \mathbf{T}])]$
HCIZ Integral	$\langle \exp(N \beta \text{Tr} [\frac{1}{N} \mathbf{T}^\top \tilde{\mathbf{A}}_2 \mathbf{T}]) \rangle_{\tilde{\mathbf{A}}}$
Layer Quality Generating Function	$\beta \mathbf{\Gamma}_{\bar{Q}^2, N \gg 1}^{IZ} := N \beta \sum_{\mu=1}^{\tilde{M}} \int_{\tilde{\lambda}_{min}}^{\tilde{\lambda}_{max}} dz R(z)$
Norm Generating Function	$G_A(\gamma) = \int_{\tilde{\lambda}_{min}^{ECS}}^{\tilde{\lambda}} R_A(z) dz$
Eigenvalue for $\mathbf{X} = \frac{1}{N} \mathbf{W}^\top \mathbf{W}$	$\lambda$ or $\lambda_i$ for $i = 1 \cdots M$
Power Law ESD Tail for $\mathbf{X}$	$\rho_{tail}(\lambda) \sim \lambda^{-\alpha}$
Effective Correlation Space ESD Tail for $\mathbf{X}$	$\rho_{tail}^{ECS}(\tilde{\lambda}), \text{Tr} \left[ \ln \prod_{j=1}^{\tilde{M}} \tilde{\lambda}_j \right] = 0$
Schatten Norm	$\ \mathbf{X}\ _\alpha^\alpha = \sum_j \lambda_j^\alpha$
ECS Tail (or Trace) Norm	$\frac{1}{\tilde{M}} \sum_i \tilde{\lambda}_i, \tilde{\lambda}_i \in \rho_{tail}^{ECS}(\tilde{\lambda})$
Spectral Norm	$\ \mathbf{X}\ _\infty = \lambda_{max}$
WeightWatcher Start of PL Tail	$\lambda_{min}^{PL} = \lambda_{min}$
Start of ECS Tail	$\lambda_{min}^{ECS} = \lambda_{min}^{ detX =1}$
ECS-PL Gap between start of tails	$\Delta \lambda_{min} := \lambda_{min}^{ECS} - \lambda_{min}^{ detX =1}$
WeightWatcher Alpha (layer) quality metric	$\alpha$
WeightWatcher AlphaHat (layer) quality metric	$\hat{\alpha} = \alpha \log_{10}(\lambda_{max})$

Table 8: Summary of various symbols used throughout the text.

Explanation	Examples	Refs
<b>Energy Landscape or Output function</b> The output of the NN given a single input data point.	$E$	Sec. 2.1 1,74,78,99
<b>Energy Difference or Student-Teacher (ST) Error</b> The difference between the output of a Student NN and its prescribed label $y$ for a single data point And as the total error for a sample of $N$ data points Or between the outputs of the Student and the Teacher NNs.	$\Delta E, \Delta \mathbf{E}$ $\Delta E := y - E_S$ $\Delta \mathbf{E} = \sum \Delta E$ $\Delta \mathbf{E} := \sum E_S - E_T$	Sec. 4, A.2 19 22 85,100,172
<b>Annealed Hamiltonian (and Potentials)</b> The Annealed Error Potential (for the Error) is defined as: The Annealed Hamiltonian (for the Error): At high-T, the relation between $H_{hT}^{an}$ and $\epsilon(R)$ is The full ST model Hamiltonian: At high-T, the ST model Hamiltonian is: The matrix-generalized ST model Hamiltonian: At high-T, the matrix-generalized ST model Hamiltonian is: The Layer Quality-Squared Hamiltonian (for the Accuracy):	$H^{an}(\mathcal{E}(R), \epsilon(R))$ $\epsilon(R) = \frac{1}{N} \mathcal{E}(R) = \langle \Delta \mathbf{E} \rangle_{\xi^N}$ $\beta H^{an} := \frac{1}{N} \ln \langle e^{-\beta \Delta \mathbf{E}} \rangle_{\xi^N}$ $H_{hT}^{an}(R) := \epsilon(R)$ $\beta H^{an}(\beta, R) := \frac{1}{2} \ln[1 + 2\beta(1 - R)]$ $\beta H_{hT}^{an}(R) := \beta(1 - R)$ $\beta H^{an}(\mathbf{R}) := \frac{1}{2} \ln[1 + 2\beta(\mathbf{I}_M - \mathbf{R})]$ $H_{hT}^{an}(\mathbf{R}) := \mathbf{I}_M - \mathbf{R}$ $\mathbf{H}_{\bar{Q}^2} := \mathbf{R}^\top \mathbf{R}$	Sec. 4.2.4, A.2 96 44, 156 49 163 96, 164 192 193 16, 205
<b>Different Average Model Errors</b> Empirical Training, Teacher, and Test Errors Empirical Generalization Gap Student-Teacher Training and Generalization Errors Neural Network (MLP) Training and Generalization Errors, the (abstract) matrix generalization of ST error	$\bar{\mathcal{E}}$ $\bar{\mathcal{E}}_{train}^{emp}, \bar{\mathcal{E}}^T \approx \bar{\mathcal{E}}_{gen}^{emp}$ $\bar{\mathcal{E}}_{gap}^{emp} := \bar{\mathcal{E}}_{train}^{emp} - \bar{\mathcal{E}}_{gen}^{emp}$ $\bar{\mathcal{E}}_{train}^{ST}, \bar{\mathcal{E}}_{gen}^{ST}$ $\bar{\mathcal{E}}_{train}^{NN}, \bar{\mathcal{E}}_{gen}^{NN}$	Sec. 4.2.5 75,76,82 83 159,162
<b>Average Training and/or Generalization Error</b> In the AA and at High-T, these are the same, and are just the Thermal Average of $\epsilon(R)$ For the ST model, we always assume AA and High-T Likewise, when generalizing $\bar{\mathcal{E}}_{gen}^{ST}$ to matrices,	$\bar{\mathcal{E}}_{train}, \bar{\mathcal{E}}_{gen}$ $\bar{\mathcal{E}}_{train}^{an, hT} = \bar{\mathcal{E}}_{gen}^{an, hT} = \langle \epsilon(R) \rangle_s^\beta$ $\bar{\mathcal{E}}_{gen}^{ST} = \bar{\mathcal{E}}_{gen}^{an, hT}$ $\bar{\mathcal{E}}_{gen}^{ST} \rightarrow \bar{\mathcal{E}}_{gen}^{NN} = \bar{\mathcal{E}}_{gen}^{an, hT}$	Sec. 4.2.5 ??,??
<b>Layer Qualities</b> For the ST Perceptron, $\bar{Q}^{ST}$ is the generalization accuracy in terms of the Self-Overlap $\eta(R)$ In the AA, and at high-T, $\langle \epsilon(R) \rangle_s^\beta = 1 - R$ For an MLP / NN, we approximate the total accuracy as a product of layer qualities $\bar{Q}$ (in the AA, at high-T) For a matrix, the Layer Quality-Squared $\bar{Q}^2$ is an HCIZ integral We approximate $\bar{Q}$ using the quality squared	$\bar{Q}, \bar{Q}^2$ $\bar{Q}^{ST} := 1 - \bar{\mathcal{E}}_{gen}^{ST} = 1 - \langle \epsilon(R) \rangle_s^\beta$ $\bar{Q}^{ST} := \langle \eta(R) \rangle_s^\beta$ $\bar{Q}^{ST} := 1 - \bar{\mathcal{E}}_{gen}^{an, ht} = \langle R \rangle_s^\beta$ $\bar{Q}^{NN} := \prod \bar{Q}_L^{NN}$ $\bar{Q}^2 := \langle \mathbf{R}^\top \mathbf{R} \rangle_s^\beta$ $\bar{Q} := \sqrt{\bar{Q}^2} \approx Q_L^{NN}$	Sec. 5.2.1, A.3 98 7 109 9,110

Table 9: Summary of types of “Energies,” with simplified examples of the notation, and references to definitions. This is a guide to understanding how the various Energies are defined and used. See the text for exact definitions, dependent variables, etc.

## A.2 Summary of the Statistical Mechanics of Generalization (SMOG)

In this section, we derive the Annealed Hamiltonian for two variants of the ST model, in the high-T limit: in Appendix A.2.1, we derive an expression for  $H^{an}(R)$  for the ST Perceptron model, when the students and teachers are modeled as  $N$ -vectors  $\mathbf{w}$  (as in [?]); and in Appendix A.2.2, we derive an expression for  $H^{an}(\mathbf{R})$  for Matrix-Generalized case, i.e., when the students and teachers are modeled as  $N \times M$  matrices  $\mathbf{W}$  (as our SETOL requires). From these, we will obtain expressions for the Average ST Model Generalization Accuracy  $\bar{\mathcal{E}}_{gen}^{ST}$  and the Average NN Model Generalization Accuracy  $\bar{\mathcal{E}}_{gen}^{NN}$ , as well as for the corresponding data-averaged errors. Although the functional form for these quantities will be the same for the vector case and the matrix case, there are several important differences in the derivation of  $H^{an}(\mathbf{R})$ , most notably having to do with a normalization for the weight matrix.

### A.2.1 Annealed Hamiltonian $H^{an}(R)$ when Student and Teachers are Vectors

In this section, we derive an expression for the Annealed Hamiltonian  $H^{an}(R)$ , in the AA and the high-T approximation, when student and teachers are modeled as  $N$ -vectors. From this, we obtain an expression for the data-averaged ST error  $\epsilon(R)$ , which is the same as the expression given in Eqn. 96.

The procedure starts by computing the associated quenched average of the Free Energy, defined for the model error as

$$\begin{aligned} \langle -\beta F \rangle_{\xi^N} &:= \langle \ln Z \rangle_{\xi^N} \\ &= \left\langle \int d\mu(\mathbf{s}) e^{-\beta \Delta \mathbf{E}_{\ell_2}(\mathbf{s}, \mathbf{t}, \xi^n)} \right\rangle_{\xi^N} \\ &= \frac{1}{N} \int d\mu(\xi^N) \ln \int d\mu(\mathbf{s}) e^{-\beta \Delta \mathbf{E}_{\ell_2}(\mathbf{s}, \mathbf{t}, \xi^n)}, \end{aligned} \quad (153)$$

where  $d\mu(\xi^N) := \prod_{i=1}^N d\xi_i P(\xi^N)$  (see Eqn. ??) and where the data-dependent ST error,  $\Delta \mathbf{E}_{\ell_2}(\mathbf{s}, \mathbf{t}, \xi^n)$ , is defined in Eqn. 85, with an  $\mathcal{L} = \ell_2$  loss. If we apply the AA (see Eqn. 41) to Eqn. 153, then we obtain

$$\langle -\beta F \rangle_{\xi^N} \simeq \ln \frac{1}{N} \int d\mu(\xi^N) \int d\mu(\mathbf{s}) e^{-\beta \Delta \mathbf{E}_{\ell_2}(\mathbf{s}, \mathbf{t}, \xi^n)}. \quad (154)$$

Notice that we have interchanged the logarithm ( $\ln$ ) and the data average (the “disorder average”)  $\langle \dots \rangle_{\xi^N}$  over the data; this is the essence of the AA, as it lets the disorder fluctuate rather than forcing the system to be quenched to the data. We will now switch the order of integration in Eqn. 154, giving

$$\langle -\beta F \rangle_{\xi^N} \simeq \ln \int d\mu(\mathbf{s}) \frac{1}{N} \int d\mu(\xi^N) e^{-\beta \Delta \mathbf{E}_{\ell_2}(\mathbf{s}, \mathbf{t}, \xi^n)}. \quad (155)$$

We now recall the definition of the Annealed Hamiltonian,  $H^{an}(R)$  (see Eqn. 44 in Section 4.2, which is analogous to Eqn. (2.31) of [?]):

$$\beta H^{an}(R) = \beta H^{an}(\beta, \mathbf{s}, \mathbf{t}) := -\frac{1}{N} \ln \int d\mu(\xi^N) e^{-\beta \Delta \mathbf{E}_{\ell_2}(\mathbf{s}, \mathbf{t}, \xi^n)}. \quad (156)$$

where we have denoted the Hamiltonian as  $H^{an}(\beta, \mathbf{s}, \mathbf{t})$  to indicate the explicit dependence on  $\beta$ , and we have added  $\beta$  to the R.H.S. because the L.H.S. is unitless. Using this definition, we can

express the Annealed Partition Function,  $Z_N^{an}$ , in the AA in terms of the Annealed Hamiltonian  $H^{an}(R)$  (as in Eqn. 46 in Section 4.2, and as in Eqn. (2.31) of [?]):

$$Z_N^{an} := N\langle Z \rangle_{\xi_N} = \int d\mu(\mathbf{s}) e^{-N\beta H^{an}(\beta, \mathbf{s}, \mathbf{t})}. \quad (157)$$

Following Section 4.2.5, we can write the *Average Model Training Error*  $\bar{\mathcal{E}}_{train}^{ST}$ , in the AA, in terms of Annealed Partition Function,  $Z_N^{an}$ :

$$\bar{\mathcal{E}}_{train}^{ST} := \frac{1}{N} \frac{\partial}{\partial \beta} Z_N^{an} \quad (158)$$

This now lets us write the Average Model Training Error  $\bar{\mathcal{E}}_{train}^{ST}$  in the AA in terms of the Annealed Hamiltonian  $H^{an}(R)$ :

$$\begin{aligned} \bar{\mathcal{E}}_{train}^{ST} &= \frac{1}{Z_N^{an}} \int d\mu(\mathbf{s}) \frac{\partial \beta H^{an}}{\partial \beta} e^{-N\beta H^{an}(\beta, \mathbf{s}, \mathbf{t})} \\ &= \frac{1}{Z_N^{an}} \int d\mu(\mathbf{s}) \langle \Delta \mathbf{E}_{\ell_2}(\mathbf{s}, \mathbf{t}, \boldsymbol{\xi}) \rangle_{\xi^{train}} e^{-N\beta H^{an}(\beta, \mathbf{s}, \mathbf{t})}, \end{aligned} \quad (159)$$

where  $\langle \mathbf{E}_{\ell_2}(\mathbf{s}, \mathbf{t}, \boldsymbol{\xi}) \rangle_{\xi^{train}}^\beta$  is a Thermal Average but defined over the specific ST error in the AA for the chosen set of training data  $\xi^{train}$ , and is denoted by

$$\langle \mathbf{E}_{\ell_2}(\mathbf{s}, \mathbf{t}, \boldsymbol{\xi}) \rangle_{\xi^{train}}^\beta := \frac{\partial \beta H^{an}}{\partial \beta}. \quad (160)$$

This can be seen from

$$\begin{aligned} \frac{\partial \beta H^{an}}{\partial \beta} &= \frac{\partial}{\partial \beta} \left( -\frac{1}{N} \ln \int d\mu(\boldsymbol{\xi}^n) e^{-\beta \Delta \mathbf{E}_{\ell_2}(\mathbf{w}, \boldsymbol{\xi}^n)} \right) \\ &= -\frac{1}{N} \frac{\partial}{\partial \beta} \frac{1}{N} \ln \int d\mu(\boldsymbol{\xi}^n) e^{-\beta \Delta \mathbf{E}_{\ell_2}(\mathbf{w}, \boldsymbol{\xi}^n)} \\ &= -\frac{1}{N} \left( \int d\mu(\boldsymbol{\xi}^n) e^{-\beta \Delta \mathbf{E}_{\ell_2}(\mathbf{w}, \boldsymbol{\xi}^n)} \right)^{-1} \frac{\partial}{\partial \beta} \int d\mu(\boldsymbol{\xi}^n) e^{-\beta \Delta \mathbf{E}_{\ell_2}(\mathbf{w}, \boldsymbol{\xi}^n)} \\ &= -\frac{1}{N} \left( \int d\mu(\boldsymbol{\xi}^n) e^{-\beta \Delta \mathbf{E}_{\ell_2}(\mathbf{w}, \boldsymbol{\xi}^n)} \right)^{-1} \int d\mu(\boldsymbol{\xi}^n) (-\Delta \mathbf{E}_{\ell_2}(\mathbf{w}, \boldsymbol{\xi}^n)) e^{-\beta \Delta \mathbf{E}_{\ell_2}(\mathbf{w}, \boldsymbol{\xi}^n)} \end{aligned} \quad (161)$$

This is analogous to defining the average error  $\Delta \mathbf{E}_{\ell_2}(\mathbf{w}, \boldsymbol{\xi}^n)$  as a Thermal Average, but as one over the data  $\boldsymbol{\xi}^n$  instead of the weights.

We can also write the Model Generalization Accuracy  $\bar{\mathcal{E}}_{gen}^{ST}$  as Boltzmann-weighted average of  $\epsilon(R)$ , weighted by  $H^{an}(\beta, S; T)$  (as in (2.30) in [?]), as:

$$\bar{\mathcal{E}}_{gen}^{ST} = \frac{1}{Z_N^{an}} \int d\mu(\mathbf{s}) \epsilon(R) e^{-N\beta H^{an}(\beta, \mathbf{s}, \mathbf{t})}, \quad (162)$$

where  $\epsilon(R) = \epsilon(\mathbf{s})$  is the average ST error, for a fixed Teacher T, averaged over *all* possible data inputs, i.e., not just over the specific training data. (Note that we have dropped the subscript *train* on  $\boldsymbol{\xi}$  since it is clear from the context.)

In the high-T (small  $\beta$ ) limit, the two model errors become formally equivalent (i.e.,  $\bar{\mathcal{E}}_{train}^{ST} = \bar{\mathcal{E}}_{gen}^{ST}$  as  $T \rightarrow \infty$ ). To show this, consider the Annealed Hamiltonian  $H^{an}(R)$ , for the Linear Perceptron with the  $\ell_2$  loss. As shown in Eqn. (C6) of [?], this takes a simple analytic form—in the large- $N$  limit—in terms of the ST overlap  $R$ :

$$H^{an}(R) = \frac{1}{2} \ln [1 + 2\beta(1 - R)]. \quad (163)$$



Eqn. 163 holds in the AA, but not in the High-T limit. If we evaluate  $\frac{\partial H^{an}}{\partial \beta}$  in the High-T (small  $\beta$ ) limit, then we can use the approximation ( $\ln[1+x] \simeq x + \dots$ ) to obtain the High-T approximation:

$$\beta H^{an}(R) \simeq \beta H_{hT}^{an}(R) := \beta(1 - R), \quad \beta \text{ small.} \quad (164)$$

where we now see that  $H_{hT}^{an}(R)$  no longer explicitly depends on  $\beta$ . By Eqn. 50, this gives

$$\epsilon(R) = \langle \mathbf{E}_{\ell_2}(\mathbf{s}, \mathbf{t}, \boldsymbol{\xi}) \rangle_{\boldsymbol{\xi}^N} \simeq 1 - R \text{ as } N \rightarrow \infty, \quad (165)$$

which we recognize as the same as the data-averaged ST error  $\epsilon(R)$  in Eqn. 96.

### A.2.2 Annealed Hamiltonian $H^{an}(\mathbf{R})$ for the Single Layer Matrix-Generalized ST Error

In this section, we derive an expression for our matrix generalization of the Annealed Hamiltonian of the Linear Perceptron, in the AA and the high-T approximation, when student and teachers are modeled as  $N \times M$  matrices  $\mathbf{W}$ , i.e.,  $H^{an}(R) \rightarrow H^{an}(\mathbf{R})$ . See Eqn. ??, which has the same form as Eqn. 156 for the vector case. From this, we obtain an expression for the data-averaged ST error  $\epsilon(R)$ , again when the student and teachers are modeled as matrices. There is a subtle normalization issue here, about which we need to be careful. However, when we normalize appropriately, we will obtain an expression for “data-averaged ST error” (i.e., Annealed Error Potential)  $\epsilon(R)$  that is of the same form as we obtained in the vector case (as given in Eqn. 165 and Eqn. 96). The difference will be that in the vector case we take  $\mathbf{R} = \frac{1}{N} \mathbf{s}^\top \mathbf{t}$ , while in the matrix case we take  $\mathbf{R} = \frac{1}{N} \mathbf{S}^\top \mathbf{T}$ .

We will need to evaluate an average over the  $N$  random  $M$ -dimensional training data vectors  $\boldsymbol{\xi}$ , which are i.i.d Gaussian with 0 mean and  $\sigma^2$  variance:

$$\|\boldsymbol{\xi}\|^2 := \frac{1}{N} \sum_{\mu=1}^N \boldsymbol{\xi}_\mu \boldsymbol{\xi}_\mu^\top = \sigma^2 \mathbf{I}_M, \quad (166)$$

where  $\boldsymbol{\xi}_\mu$  is a vector of length  $M$ , and  $\mathbf{I}_M$  is an  $M \times M$  identity matrix. The expected value of the squared norm is:

$$\mathbb{E}[\|\boldsymbol{\xi}_\mu\|^2] = M\sigma^2. \quad (167)$$

If we let  $\sigma^2 \sim \frac{1}{M}$ , then  $\mathbb{E}[\|\boldsymbol{\xi}_\mu\|^2] = 1$ , i.e., the data vectors can be normalized to 1. Let the probability distribution over the  $N$  data vectors be

$$\begin{aligned} P(\boldsymbol{\xi}^N) &= \prod_{\mu=1}^N \left( \frac{1}{\sqrt{(2\pi\sigma^2)^M}} \right) e^{-\frac{\|\boldsymbol{\xi}_\mu\|^2}{2\sigma^2}} \\ &= \left( \frac{1}{\sqrt{(2\pi\sigma^2)^M}} \right)^N \exp \left[ -\frac{M}{2} \sum_{\mu=1}^N \|\boldsymbol{\xi}_\mu\|^2 \right] \\ &= \mathcal{N} \exp \left[ -\frac{M}{2} \sum_{\mu=1}^N \|\boldsymbol{\xi}_\mu\|^2 \right], \end{aligned} \quad (168)$$

where  $M = N_f$  is the number of features in the data, where the normalization  $\mathcal{N}$  is

$$\mathcal{N} := \left( \frac{1}{\sqrt{(2\pi\sigma^2)^M}} \right)^N = \left( \frac{M}{2\pi} \right)^{MN/2}. \quad (169)$$

2874 **The Total Data Sample Error ( $\Delta \mathbf{E}_{\ell_2}(\mathbf{S}, \mathbf{T})$ ) and the Matrix Normalization** First, let  
 2875 us express the matrix-generalized Total Data Sample Error,  $\Delta \mathbf{E}_{\ell_2}(\mathbf{S}, \mathbf{T})$ , for a single layer, in  
 2876 operator form (for each training example)

$$\frac{1}{N} \Delta \mathbf{E}_{\ell_2}(\mathbf{S}, \mathbf{T}) := \text{Tr} \left[ \mathbf{I}_M - \frac{1}{N} \mathbf{S}^\top \mathbf{T} \right] = M - \text{Tr}[\mathbf{R}] \quad (170)$$

2877 where  $\mathbf{I}_M$  is a diagonal matrix of dimension  $M$ . Notice that  $\Delta \mathbf{E}_{\ell_2}(\mathbf{S}, \mathbf{T})$  scales as  $N \times M$ , the total  
 2878 number of parameters in the system. Also, importantly, in this representation, the normalization  
 2879 on becomes

$$\frac{1}{N} \text{Tr}[\mathbf{S}^\top \mathbf{T}] = \text{Tr}[\mathbf{R}] = M \quad (171)$$

2880 so that when all  $M$  elements overlap, then the error is zero. We can define the data-dependent  
 2881 form (i.e., in the basis of the data  $\boldsymbol{\xi}$ ) as

$$\begin{aligned} \Delta \mathbf{E}_{\ell_2}(\mathbf{S}, \mathbf{T}, \boldsymbol{\xi}) &:= \sum_{\mu=1}^N (\boldsymbol{\xi}^\mu)^\top \left( \mathbf{I}_M - \frac{1}{N} \mathbf{S}^\top \mathbf{T} \right) \boldsymbol{\xi}^\mu \\ &= \sum_{\mu=1}^N \sum_{i,j=1}^M \boldsymbol{\xi}_i^\mu \left( \delta_{ij} - \frac{1}{N} [\mathbf{S}^\top \mathbf{T}]_{ij} \right) \boldsymbol{\xi}_j^\mu. \end{aligned} \quad (172)$$

2882 **The Annealed Hamiltonian (per-parameter,  $h^{an}(\mathbf{R})$ )** The definition of the Annealed  
 2883 Hamiltonian,  $H^{an}(\mathbf{R})$ , for the matrix-generalized case must be extended to account for the  $M$   
 2884 parameters per training example, and is now given as

$$H^{an}(\mathbf{R}) := M h^{an}(\mathbf{R}) \quad (173)$$

2885 where the Annealed Hamiltonian per-parameter,  $h^{an}(\mathbf{R})$ , is obtained from Eqn. 156 as

$$\begin{aligned} \beta h^{an}(\mathbf{R}) &:= -\frac{1}{N} \ln \int \mathcal{D}\boldsymbol{\xi}^N e^{-\beta \Delta \mathbf{E}_{\ell_2}(\mathbf{S}, \mathbf{T})} P(\boldsymbol{\xi}^N) \\ &= -\frac{1}{N} \ln \mathbb{I}_H, \end{aligned} \quad (174)$$

2886 where

$$\mathbb{I}_H := \int \mathcal{D}\boldsymbol{\xi}^N e^{-\beta \Delta \mathbf{E}_{\ell_2}(\mathbf{S}, \mathbf{T})} P(\boldsymbol{\xi}^N). \quad (175)$$

2887 That is,  $h^{an}(\mathbf{R})$  represents the Energy or Error each of the  $M$  parameters contributes (integrated  
 2888 over the  $N$  training examples  $\boldsymbol{\xi}^N$ ).

2889 The goal will be to derive the high-Temperature Annealed Hamiltonian,  $H_{hT}^{an}(\mathbf{R})$ , which is now

$$H_{hT}^{an}(\mathbf{R}) := M H_{hT}^{an}(\mathbf{R}) \quad (176)$$

2890 If we were to evaluate the trace of  $H_{hT}^{an}(\mathbf{R})$ , then we could also define a matrix-generalized  
 2891 Annealed Error Potential,

$$\epsilon(\mathbf{R}) := \text{Tr}[H_{hT}^{an}(\mathbf{R})], \quad (177)$$

2892 which would be like a mean-field potential, now summed over all the  $M$  parameters.

2893 To evaluate the integral, notice that  $\mathbb{I}_H$  is really just an average over i.i.d. data, and so it is  
 2894 just product over  $N$  independent terms

$$\mathbb{I}_H := \int \mathcal{D}\boldsymbol{\xi}^N e^{-\beta \Delta \mathbf{E}_{\ell_2}(\mathbf{S}, \mathbf{T})} P(\boldsymbol{\xi}^N) \rightarrow \left[ \int \mathcal{D}\boldsymbol{\xi} [\dots] \right]^N, \quad (178)$$

2895 as in Eqn. 182 below. Moreover, when taking  $\ln \mathbb{I}_H$ , the  $N$  term pulls down and become a  
 2896 prefactor

$$-\ln \mathbb{I}_H = -\ln \left[ \int \mathcal{D}\boldsymbol{\xi} [\dots] \right]^N = -N \ln \left[ \int \mathcal{D}\boldsymbol{\xi} [\dots] \right]. \quad (179)$$

2897 Thus, as with the vector case,  $H^{an}(\mathbf{R})$  is like a mean-field average over the data  $\boldsymbol{\xi}$ , independent of  
 2898 the sample size  $N$ . Also, since the final result must scale as  $N \times M$ , the integral should scale as  
 2899  $M$ , i.e.,  $[\int \mathcal{D}\boldsymbol{\xi} [\dots]] \sim M$ .

2900 If we substitute  $\Delta \mathbf{E}_{\ell_2}(\mathbf{S}, \mathbf{T}, \boldsymbol{\xi})$ , Eqn. 172, into the integral  $\mathbb{I}_H$ , Eqn. 175, then we obtain

$$\mathbb{I}_H = \int \mathcal{D}\boldsymbol{\xi}^N \exp \left( -\beta \sum_{\mu=1}^N (\boldsymbol{\xi}^\mu)^\top \left( \mathbf{I}_M - \frac{1}{N} \mathbf{S}^\top \mathbf{T} \right) \boldsymbol{\xi}^\mu \right) P(\boldsymbol{\xi}^N) \quad (180)$$

$$\begin{aligned} &= \int \mathcal{D}\boldsymbol{\xi}^N \exp \left( -\beta \sum_{\mu=1}^N (\boldsymbol{\xi}^\mu)^\top \left( \mathbf{I}_M - \frac{1}{N} \mathbf{S}^\top \mathbf{T} \right) \boldsymbol{\xi}^\mu \right) \mathcal{N} \exp \left( -\sum_{\mu=1}^N \frac{\|\boldsymbol{\xi}^\mu\|^2}{2\sigma^2} \right) \\ &= \mathcal{N} \int \mathcal{D}\boldsymbol{\xi}^N \exp \left( -\beta \sum_{\mu=1}^N (\boldsymbol{\xi}^\mu)^\top \left( \mathbf{I}_M - \frac{1}{N} \mathbf{S}^\top \mathbf{T} \right) (\boldsymbol{\xi}^\mu) - \sum_{\mu=1}^N \frac{\|\boldsymbol{\xi}^\mu\|^2}{2\sigma^2} \right) \\ &= \mathcal{N} \int \mathcal{D}\boldsymbol{\xi}^N \exp \left( -\frac{1}{2\sigma^2} \sum_{\mu=1}^N 2\beta\sigma^2 (\boldsymbol{\xi}^\mu)^\top \left( \mathbf{I}_M - \frac{1}{N} \mathbf{S}^\top \mathbf{T} \right) (\boldsymbol{\xi}^\mu) + \|\boldsymbol{\xi}^\mu\|^2 \right) \\ &= \mathcal{N} \int \mathcal{D}\boldsymbol{\xi}^N \exp \left( -\frac{1}{2\sigma^2} \sum_{\mu=1}^N (\boldsymbol{\xi}^\mu)^\top [2\beta\sigma^2 (\mathbf{I}_M - \frac{1}{N} \mathbf{S}^\top \mathbf{T}) + \mathbf{I}_M] (\boldsymbol{\xi}^\mu) \right). \end{aligned} \quad (181)$$

2901 By combining the exponents, we obtain

$$\begin{aligned} \mathbb{I}_H &= \mathcal{N} \int \mathcal{D}\boldsymbol{\xi}^N \exp \left[ -\frac{1}{2\sigma^2} \sum_{\mu=1}^N (\boldsymbol{\xi}^\mu)^\top (\mathbf{M}) \boldsymbol{\xi}^\mu \right] \\ &= \mathcal{N} \int \mathcal{D}\boldsymbol{\xi} \exp \left[ -\frac{1}{2\sigma^2} (\boldsymbol{\xi})^\top (\mathbf{M}) \boldsymbol{\xi} \right]^N, \end{aligned} \quad (182)$$

2902 where  $\mathbf{M} = 2\beta\sigma^2 (\mathbf{I}_M - \frac{1}{N} \mathbf{S}^\top \mathbf{T}) + \mathbf{I}_M$  is an  $M \times M$  matrix.

2903 We now use the familiar property of multi-variant Gaussian integrals,

$$\int d\mathbf{x} e^{-\frac{1}{2\sigma^2} (\mathbf{x})^\top \mathbf{M} (\mathbf{x})} = (2\pi\sigma^2)^{m/2} \frac{1}{\sqrt{\det(\mathbf{M})}} \quad (183)$$

2904 where  $\mathbf{x}$  is an  $m$ -dim vector (with zero mean), and  $\mathbf{M}$  is a square postive-definite matrix, and  
 2905  $\det(\mathbf{M})$  is the determinant of  $\mathbf{M}$ . Using Eqn. 183, we can rewrite  $\mathbb{I}_H$  in Eqn. 182 as

$$\mathbb{I}_H = \mathcal{N} \left[ \frac{(2\pi\sigma^2)^{M/2}}{\sqrt{\det(\mathbf{M})}} \right]^N \quad (184)$$

$$\begin{aligned} &= \mathcal{N} (2\pi\sigma^2)^{NM/2} \left[ \sqrt{\det \left( 2\beta\sigma^2 (\mathbf{I}_M - \frac{1}{N} \mathbf{S}^\top \mathbf{T}) + \mathbf{I}_M \right)} \right]^{-N} \\ &= \left( \frac{1}{2\pi\sigma^2} \right)^{MN/2} (2\pi\sigma^2)^{M/2} \left[ \sqrt{\det \left( \mathbf{I}_M + 2\beta\sigma^2 (\mathbf{I}_M - \frac{1}{N} \mathbf{S}^\top \mathbf{T}) \right)} \right]^{-N}, \end{aligned} \quad (185)$$

where Eqn. 185 follows by inserting  $\mathcal{N}$  from Eqn. 169. We can now identify  $\sigma^2 = \frac{1}{M}$  to obtain

$$\begin{aligned}\mathbb{I}_H &= \left[ \sqrt{\det(\mathbf{I}_M + 2\beta\sigma^2(\mathbf{I}_M - \frac{1}{N}\mathbf{S}^\top\mathbf{T}))} \right]^{-N} \\ &= \left[ \sqrt{\det(\mathbf{I}_M + \frac{2\beta}{MN}(\mathbf{I}_M - \frac{1}{N}\mathbf{S}^\top\mathbf{T}))} \right]^{-N} \\ &= \left[ \det(\mathbf{I}_M + \frac{2\beta}{M}(\mathbf{I}_M - \frac{1}{N}\mathbf{S}^\top\mathbf{T})) \right]^{-N/2}.\end{aligned}\quad (186)$$

**Large- $N$  Approximation.** Using the expression  $\det(\mathbf{I}_M + \epsilon\mathbf{\Omega}) \approx 1 + \epsilon\text{Tr}[\mathbf{\Omega}]$ , which holds for an arbitrary matrix  $\mathbf{\Omega}$  for small  $\epsilon$ , we can evaluate the determinant in Eqn. 186 in the large- $N$  approximation, which gives

$$\mathbb{I}_H \approx \left[ 1 + \frac{2\beta}{M}(\text{Tr}[\mathbf{I}_M - \frac{1}{N}\mathbf{S}^\top\mathbf{T}]) \right]^{-N/2}.\quad (187)$$

Inserting this into Eqn. 174, we obtain

$$\beta h^{an}(\mathbf{R}) = -\frac{1}{N} \ln \left[ 1 + \frac{2\beta}{M}(\text{Tr}[\mathbf{I}_M - \frac{1}{N}\mathbf{S}^\top\mathbf{T}]) \right]^{-N/2}.\quad (188)$$

**Matrix-Generalized ST Error  $H_{hT}^{an}(\mathbf{R})$  when  $M = 1$ .** To start, observe that when  $M = 1$ , Eqn. 188 becomes

$$\begin{aligned}\beta H^{an}(\mathbf{R})|_{M=1} &= \beta h^{an}(\mathbf{R})|_{M=1} \\ &= -\frac{1}{N} \ln \left[ 1 + 2\beta(1 - \frac{1}{N}\mathbf{S}^\top\mathbf{T}) \right]^{-N/2}\end{aligned}\quad (189)$$

$$\begin{aligned}&= \frac{1}{2} \ln \left[ 1 + 2\beta(1 - \frac{1}{N}\mathbf{S}^\top\mathbf{T}) \right] \\ &= \frac{1}{2} \ln [1 + 2\beta(1 - \mathbf{R})],\end{aligned}\quad (190)$$

meaning that Eqn. 188 reduces to Eqn. 163, as desired.

**FIX THIS SECTION** This ensures the Hamiltonian scales as  $M$  so the Free Energy scales as  $N \times M$ , the number of free paramaters in the system. Notice that for the final Layer Quality-Squared Hamiltonian  $\mathbf{H}_{\tilde{Q}^2}$ , this will change.

$$\mathbb{I}_H \approx \left[ 1 + \frac{2\beta}{M} \text{Tr}[(\mathbf{I}_M - \frac{1}{N}\mathbf{S}^\top\mathbf{T})] \right]^{-N/2},\quad (191)$$

for any  $M > 1$ . Given this, it follows from Eqn. eqn:HANPP and Eqn. 174 that

$$\beta H^{an}(\mathbf{R}) = \beta M h^{an}(\mathbf{R}) = \frac{M}{2} \ln \left[ 1 + \frac{2\beta}{M}(\mathbf{I}_M - \mathbf{R}) \right],\quad (192)$$

which is of the same functional form as Eqn. 163, as desired:

$$H_{hT}^{an}(\mathbf{R}) = \mathbf{I}_M - \mathbf{R} = M - \mathbf{R}\quad (193)$$

This form of the Hamiltonian,  $H^{an}(\mathbf{R})$ , however, is not symmetric, and we will eventually want a symmetric operator or matrix. Fortunately, the high-T form,  $H_{hT}^{an}(\mathbf{R})$ , can be made symmetric, as shown below.

### 2922 A.3 Expressing the Layer Quality

2923 In this section, we obtain an approximation expression for the Layer Quality-Squared from the IZ  
2924 Free Energy for the Generalization Error, given in Eqn. 112 in Section 5.2.1.

2925 For the required Free Energy  $\beta\mathbf{F}^{IZ}$ , we will use the matrix-generalized Hamiltonian from  
2926 Eqn. 193 for the Layer Quality,  $H_{hT}^{an}(\mathbf{R}) = \mathbb{I}_M - \mathbf{R}$ , giving a Boltzmann distribution and the  
2927 corresponding Thermal Average. Expanding this out, we have

$$-\beta\mathbf{F}^{IZ} = -\ln \int d\mu(\mathbf{S}) \exp[-N\beta \text{Tr}[H_{hT}^{an}(\mathbf{R})]] \quad (194)$$

$$(195)$$

2928 We could also express  $\beta\mathbf{F}^{IZ}$  In terms of the matrix-generalized Annealed Error Potential  $\epsilon(\mathbf{R})$   
2929 (Eqn. 177), giving

$$-\beta\mathbf{F}^{IZ} = -\ln \int d\mu(\mathbf{S}) \exp[-N\beta\epsilon(\mathbf{R})] \quad (196)$$

2930 In analogy with Eqn. 49, as  $H_{hT}^{an}(\mathbf{R}) = M - \mathbf{R}$ , write

$$-\beta\mathbf{F}^{IZ} = -\ln \int d\mu(\mathbf{S}) \exp[-N\beta \text{Tr}[M - \mathbf{R}]] \quad (197)$$

2931 Using the approximation  $\text{Tr}[\mathbf{R}] \approx \sqrt{\text{Tr}[\mathbf{R}^\top \mathbf{R}]}$ , we have

$$-\beta\mathbf{F}^{IZ} \approx -\ln \int d\mu(\mathbf{S}) \exp[-N\beta(M - \sqrt{\text{Tr}[\mathbf{R}^\top \mathbf{R}]})] \quad (198)$$

$$= -\ln \int d\mu(\mathbf{S}) \exp[-NM\beta] \exp[N\beta\sqrt{\text{Tr}[\mathbf{R}^\top \mathbf{R}]}], \quad (199)$$

$$= -\ln e^{-NM\beta} \int d\mu(\mathbf{S}) \exp[N\beta\sqrt{\text{Tr}[\mathbf{R}^\top \mathbf{R}]}], \quad (200)$$

$$= -\ln e^{-NM\beta} - \ln \int d\mu(\mathbf{S}) \exp[N\beta\sqrt{\text{Tr}[\mathbf{R}^\top \mathbf{R}]}], \quad (201)$$

$$= NM\beta - \ln \int d\mu(\mathbf{S}) \exp[\beta N\sqrt{\text{Tr}[\mathbf{R}^\top \mathbf{R}]}], \quad (202)$$

2932 Notice that, as expected, the Free Energy scales  $\beta\mathbf{F}^{IZ}$  as  $N \times M$ . Since Eqn. 194 equals  
2933 Eqn. 199, we can write the Free Energy in terms of  $\text{Tr}[\mathbf{R}^\top \mathbf{R}]$ . From Eqn. 202, we can identify a  
2934 generating function ( $\Gamma_{\bar{\mathcal{Q}}}$ ) for the layer accuracy, or Quality. For example, to compute the average  
2935 Quality  $\bar{\mathcal{Q}}$ , we would use

$$\beta\Gamma_{\bar{\mathcal{Q}}}^{IZ} := \ln \int d\mu(\mathbf{S}) \exp[N\beta\sqrt{\text{Tr}[\mathbf{R}^\top \mathbf{R}]}], \quad (203)$$

2936 and to compute the average Quality (squared)  $\bar{\mathcal{Q}}^2$ , we would use

$$\beta\Gamma_{\bar{\mathcal{Q}}^2}^{IZ} := \ln \int d\mu(\mathbf{S}) \exp[N\beta \text{Tr}[\mathbf{R}^\top \mathbf{R}]]. \quad (204)$$

2937 We have recovered Eqn. 112. We can now also define the Layer Quality-Squared Hamiltonian as

$$\mathbf{H}_{\bar{\mathcal{Q}}^2} := \mathbf{R}^\top \mathbf{R} \quad (205)$$

2938 which is a symmetric operator, as desired. Consequently, we may also write

$$\beta\Gamma_{\bar{\mathcal{Q}}^2}^{IZ} := \ln \int d\mu(\mathbf{S}) \exp[N\beta \text{Tr}[\mathbf{H}_{\bar{\mathcal{Q}}^2}]]. \quad (206)$$

## A.4 Derivation of the TRACE-LOG Condition

### A.4.1 Setting up the Saddle Point Approximation (SPA)

As in Eqn. 117, we can write Eqn. 11 in terms of the  $\mathbf{A}_2$  form of the Student Correlation matrix, giving

$$\beta \mathbf{\Gamma}_{\mathcal{Q}^2}^{IZ} = \ln \int_{\mathbf{S}} d\mu(\mathbf{S}) \exp \left( N\beta \text{Tr} \left[ \frac{1}{N} \mathbf{T}^\top \mathbf{A}_2 \mathbf{T} \right] \right) \quad (207)$$

where  $d\mu(\mathbf{S})$  is the measure over all  $N \times M$  real-valued random matrices, although we really want to limit this to all  $N \times M$  real matrices that resemble the Teacher  $\mathbf{T}$ , which we clarify below.

To transform  $\beta \mathbf{\Gamma}_{\mathcal{Q}^2}^{IZ}$  into a form we can evaluate using Tanakas result, we need to change the measure from an integral over all random  $N \times M$  student weight matrices  $d\mu(\mathbf{S})$  to an integral over all  $N \times N$  student correlation matrices  $d\mu(\mathbf{A})$ , i.e.,  $d\mu(\mathbf{S}) \rightarrow d\mu(\mathbf{A})$ . To accomplish this, we can insert an integral over the Dirac Delta function

$$\mathbf{I} := \int d\mu(\mathbf{A}_1) \delta(N\mathbf{A}_1 - \mathbf{S}^\top \mathbf{S}) = \int d\mu(\mathbf{A}) \delta(N\mathbf{A}_1 - \mathbf{S}^\top \mathbf{S}). \quad (208)$$

(This is simply a resolution of the Identity.) This gives

$$\beta \mathbf{\Gamma}_{\mathcal{Q}^2}^{IZ} = \ln \int_{\mathbf{S}} d\mu(\mathbf{S}) \int_{\mathbf{A}} d\mu(\mathbf{A}) \delta(N\mathbf{A}_1 - \mathbf{S}^\top \mathbf{S}) e^{N\beta \text{Tr} \left[ \frac{1}{N} \mathbf{T}^\top \mathbf{A}_2 \mathbf{T} \right]}, \quad (209)$$

where  $d\mu(\mathbf{A}) = \mathbf{Pr}[\mathbf{A}] d\mathbf{A}$  and  $\mathbf{Pr}[\mathbf{A}]$  is the (still unspecified) probability density over the new random matrix  $\mathbf{A}$ . Let us express Eqn. 209 at large- $N$  as

$$\lim_{N \gg 1} \beta \mathbf{\Gamma}_{\mathcal{Q}^2}^{IZ} = \lim_{N \gg 1} \ln \int d\mu(\mathbf{A}) \int d\mu(\mathbf{S}) \delta(N\mathbf{A}_1 - \mathbf{S}^\top \mathbf{S}) e^{N\beta \text{Tr} \left[ \frac{1}{N} \mathbf{T}^\top \mathbf{A}_2 \mathbf{T} \right]}. \quad (210)$$

Now we assume we can first evaluate the term

$$\lim_{N \gg 1} \int d\mu(\mathbf{S}) \delta(N\mathbf{A}_1 - \mathbf{S}^\top \mathbf{S}) \quad (211)$$

at large- $N$  using a Saddle Point Approximation (SPA).

Using the relation,

$$\delta(N\mathbf{A}_1 - \mathbf{S}^\top \mathbf{S}) = \mathcal{N} \int_{\hat{\mathbf{A}}} d\mu(\hat{\mathbf{A}}) e^{iN \text{Tr}[\hat{\mathbf{A}} \mathbf{A}_1]} e^{-i \text{Tr}[\hat{\mathbf{A}} \mathbf{S}^\top \mathbf{S}]}, \quad (212)$$

where  $\hat{\mathbf{A}}$  is a  $M \times M$  auxiliary matrix, and the domain of integration  $d\mu(\hat{\mathbf{A}})$  is all  $M \times M$  real-valued matrices, and where the normalization  $\mathcal{N}_1$  is

$$\mathcal{N}_1 := \frac{1}{(2\pi)^{M(M+1)/4}}, \quad (213)$$

because  $\mathbf{A}$  is a symmetric matrix with  $M(M+1)/2$  constraints.

This is simply the matrix generalization of  $\delta(x) = \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{i\hat{x}x} d\hat{x}$ , so we can express the delta function as an exponential, giving

$$\beta \mathbf{\Gamma}_{\mathcal{Q}^2}^{IZ} = \mathcal{N}_1 \ln \int_{\mathbf{S}} d\mu(\mathbf{S}) \int_{\mathbf{A}} d\mu(\mathbf{A}_1) \int_{\hat{\mathbf{A}}} d\mu(\hat{\mathbf{A}}) e^{iN \text{Tr}[\hat{\mathbf{A}} \mathbf{A}_1]} e^{-i \text{Tr}[\hat{\mathbf{A}} \mathbf{S}^\top \mathbf{S}]} e^{N\beta \text{Tr} \left[ \frac{1}{N} \mathbf{T}^\top \mathbf{A}_2 \mathbf{T} \right]}. \quad (214)$$

Rearranging terms, we obtain

$$\beta \mathbf{\Gamma}_{\mathcal{Q}^2}^{IZ} = \ln \int_{\mathbf{A}} d\mu(\mathbf{A}) e^{N\beta \text{Tr} \left[ \frac{1}{N} \mathbf{T}^\top \mathbf{A}_2 \mathbf{T} \right]} \times \Gamma_1, \quad (215)$$



2961 where we define  $\Gamma_1$  as

$$\Gamma_1 := \Gamma_1(\mathbf{A}_1) = \mathcal{N}_1 \int_{\mathbf{S}} d\mu(\mathbf{S}) \int_{\hat{\mathbf{A}}} d\mu(\hat{\mathbf{A}}) e^{iN\text{Tr}[\hat{\mathbf{A}}\mathbf{A}_1]} e^{-i\text{Tr}[\hat{\mathbf{A}}\mathbf{S}^\top \mathbf{S}]}.$$
 (216)

2962 We can simplify the complex integral in  $\Gamma_1$  with the Wick Rotation  $i\hat{\mathbf{A}} \rightarrow \hat{\mathbf{A}}$ . We may expect  
 2963  $d\mu(\hat{\mathbf{A}})$  to be invariant to rotations in the complex plane so the Wick rotation does not introduce  
 2964 any complex prefactors. This gives

$$\Gamma_1 = \mathcal{N}_1 \int_{\mathbf{S}} d\mu(\mathbf{S}) \int_{\hat{\mathbf{A}}} d\mu(\hat{\mathbf{A}}) e^{N\text{Tr}[\hat{\mathbf{A}}\mathbf{A}_1]} e^{-\text{Tr}[\hat{\mathbf{A}}\mathbf{S}^\top \mathbf{S}]} \quad (217)$$

$$= \mathcal{N}_1 \int_{\mathbf{S}} d\mu(\mathbf{S}) \int_{\hat{\mathbf{A}}} d\mu(\hat{\mathbf{A}}) e^{N\text{Tr}[\hat{\mathbf{A}}\mathbf{A}_1]} e^{-\text{Tr}[\mathbf{S}\hat{\mathbf{A}}\mathbf{S}^\top]}, \quad (218)$$

2965 where the second line follows since the trace is invariant under cyclic permutations (i.e.,  $\text{Tr}[ABC] =$   
 2966  $\text{Tr}[BCA] = \text{Tr}[CAB]$ ). Swapping the order of the integrals yields

$$\Gamma_1 = \Gamma_1(\hat{\mathbf{A}}) = \mathcal{N}_1 \int_{\hat{\mathbf{A}}} d\mu(\hat{\mathbf{A}}) e^{N\text{Tr}[\hat{\mathbf{A}}\mathbf{A}_1]} \times \Gamma_2, \quad (219)$$

2967 where we define  $\Gamma_2$  as

$$\Gamma_2 := \Gamma_2(\hat{\mathbf{A}}) = \int_{\mathbf{S}} d\mu(\mathbf{S}) e^{-\text{Tr}[\mathbf{S}\hat{\mathbf{A}}\mathbf{S}^\top]}.$$

2968 To evaluate  $\Gamma_2$ , we will make several mathematically convenient approximations. (These will  
 2969 yield an approximate expression which can be verified empirically.) We first assume for the  
 2970 purpose of changing measure that the (data) columns of  $\mathbf{S}$  are statistically independent, so that  
 2971 the measure  $d\mu(\mathbf{S})$  factors into  $N$  gaussian distributions

$$d\mu(\mathbf{S}) = \prod_{\mu=1}^N d\mu(\mathbf{s}_\mu) = \prod_{\mu=1}^N d\mathbf{s}_\mu, \quad (220)$$

2972 where  $\mathbf{s}_\mu$  is an  $M$ -dimensional vector. The singular values of  $\mathbf{S}$  are invariant to randomly permuting  
 2973 the columns or rows, so the resulting ESD does not change. This is very different from permuting  
 2974  $\mathbf{S}$  element-wise, which will make the resulting ESD Marchenko Pastur (MP).

2975 Using Eqn. 220,  $\Gamma_2$  reduces to a simple Gaussian integral, which can be evaluated as a product  
 2976 of  $N$  Gaussian integrals (over the  $M \times M$  matrix  $\hat{\mathbf{A}}$ )

$$\Gamma_2 = \left[ \int_{\mathbf{s}} d\mathbf{s} e^{-\frac{1}{\sigma^2} \mathbf{s} \hat{\mathbf{A}} \mathbf{s}^\top} \right]^N \quad (221)$$

$$= \left[ \mathcal{N}_2 \det(\hat{\mathbf{A}})^{-1/2} \right]^N, \quad (222)$$

2977 where the normalization  $\mathcal{N}_2$

$$\mathcal{N}_2 := (\pi\sigma^2)^{M/2}, \quad (223)$$

2978 where  $\sigma^2 = \mathbf{s}^\top \mathbf{s} = 1/M$

2979 For any square, non-singular matrix  $\hat{\mathbf{A}}$ ,  $\text{Tr}[\ln \hat{\mathbf{A}}] = \ln \det(\hat{\mathbf{A}})$ , so it follows from Eqn. 221 that

$$\begin{aligned} \ln \Gamma_2 &= N \ln \mathcal{N}_2 \left[ (\det \hat{\mathbf{A}})^{-1/2} \right] \\ &= N \ln \mathcal{N}_2 - \frac{N}{2} \text{Tr}[\ln \hat{\mathbf{A}}], \end{aligned} \quad (224)$$

2980 so that

$$\Gamma_2 = \mathcal{N}_2^N e^{-\frac{N}{2} \text{Tr}[\ln \hat{\mathbf{A}}]} \quad (225)$$

2981 Substituting Eqn. 225 into Eqn. 219, we can write  $\Gamma_1$  as

$$\Gamma_1(\hat{\mathbf{A}}) = C_{\Gamma_1} \int_{\hat{\mathbf{A}}} d\mu(\hat{\mathbf{A}}) e^{N \text{Tr}[\hat{\mathbf{A}} \mathbf{A}_1]} e^{-\frac{N}{2} \text{Tr}[\ln \hat{\mathbf{A}}]}, \quad (226)$$

2982 where

$$C_{\Gamma_1} := \mathcal{N}_1 e^{N^2}. \quad (227)$$

2983 We now can evaluate the integral in Eqn. 219 over the Lagrange Multiplier  $\hat{\mathbf{A}}$  (i.e.,  $\int_{\hat{\mathbf{A}}}$ ). If we  
2984 call this  $\Gamma_1(\hat{\mathbf{A}})$ , then (following Tanaka [?]) we can define the *Rate Function*  $I(\hat{\mathbf{A}}, \mathbf{A}_1)$  such that

$$\Gamma_1(\hat{\mathbf{A}}) = \int_{\hat{\mathbf{A}}} d\mu(\hat{\mathbf{A}}) e^{-N I(\hat{\mathbf{A}}, \mathbf{A}_1)}, \quad (228)$$

2985 where

$$I(\hat{\mathbf{A}}, \mathbf{A}_1) = -\text{Tr}[\hat{\mathbf{A}} \mathbf{A}_1] + \frac{1}{2} \text{Tr}[\ln \hat{\mathbf{A}}]. \quad (229)$$

2986 We can formally evaluate the integral in Eqn. 228 in the large- $N$  limit using a Saddle Point  
2987 Approximation (SPA) (see Section 4.2, Eqn. 67), as

$$\Gamma_1(\hat{\mathbf{A}}) \rightarrow \sqrt{\frac{(2\pi)^{N/2}}{N \|I\|}} e^{-N I^*(\hat{\mathbf{A}}, \mathbf{A}_1)}, \quad (230)$$

2988 where  $I^*(\hat{\mathbf{A}}, \mathbf{A}_1)$  is the maximum value, obtained using

$$I^*(\hat{\mathbf{A}}, \mathbf{A}_1) := \lim_{N \gg 1} I(\hat{\mathbf{A}}, \mathbf{A}_1) = \sup_{\hat{\mathbf{A}}} \left[ -\text{Tr}[\hat{\mathbf{A}} \mathbf{A}_1] + \frac{1}{2} \text{Tr}[\ln \hat{\mathbf{A}}] \right], \quad (231)$$

2989 where  $I$  at the SPA is defined as

$$I := \frac{\partial}{\partial \hat{\mathbf{A}}} I(\hat{\mathbf{A}}, \mathbf{A}_1) = -\mathbf{A}_1 + \frac{1}{2\hat{\mathbf{A}}} = 0 \quad (232)$$

2990 and  $I$  is defined as

$$I = \frac{\partial^2}{\partial \hat{\mathbf{A}}^2} I(\hat{\mathbf{A}}, \mathbf{A}_1) = -\frac{1}{2} \left( \frac{1}{2} \hat{\mathbf{A}}^{-1} \right) \otimes \left( \frac{1}{2} \hat{\mathbf{A}}^{-1} \right) = -\frac{1}{8} \mathbf{A}_1 \otimes \mathbf{A}_1 \quad (233)$$

2991 where  $\otimes$  is the Kronecker product, and  $\mathbf{A}_1 \otimes \mathbf{A}_1$  is the Hessian of  $\mathbf{A}_1$ .

2992 Solving the SPA equation, we find that the auxiliary matrix is  $\hat{\mathbf{A}} = \frac{1}{2} \mathbf{A}_1^{-1}$  and the prefactor  
2993 (Hessian) is given as  $\det\left(-\frac{1}{8} \mathbf{A}_1 \otimes \mathbf{A}_1\right) = \left(-\frac{1}{8}\right)^{M^2} (\det(\mathbf{A}_1))^M$ .

2994 [ double check prefactors.]

2995 Substituting for  $\hat{\mathbf{A}}$  into Rate Function (Eqn. 229),  $I$  becomes

$$\begin{aligned} I^*(\hat{\mathbf{A}}, \mathbf{A}_1) &= -\text{Tr}[\mathbb{I}_M] + \frac{1}{2} \text{Tr}[\ln \mathbf{A}_1] \\ &= -M + \frac{1}{2} \text{Tr}[\ln \mathbf{A}_1]. \end{aligned} \quad (234)$$

2996 In order for this result to be physically meaningful, we need that if  $I^*(\hat{\mathbf{A}}, \mathbf{A}_1)$  grows, then it  
2997 must grow slower than  $N$ , and, more importantly, that  $\det(\mathbf{A})$  be non-zero. Importantly, When  
2998  $\det(\mathbf{A}) = 1$  exactly, however, then  $\Gamma_1$  becomes a constant, and this simplifies things considerably!

#### 2999 A.4.2 Casting the Generating Function ( $\beta\Gamma_{\mathcal{Q}^2}^{IZ}$ ) as an HCIZ Integral

3000 In this section, we express the Generating Function  $\beta\Gamma_{\mathcal{Q}^2}^{IZ}$ , given in Eqn. 117 (equivalently, in  
3001 Eqn. 112), as an HCIZ Integral, as given in Eqn. 121.

3002 Inserting  $I^*(\hat{\mathbf{A}}, \mathbf{A})$  from Eqn. 234 into  $\beta\Gamma_{\mathcal{Q}^2}^{IZ}$ , we obtain

$$\begin{aligned}\beta\Gamma_{\mathcal{Q}^2}^{IZ} &= \ln \left[ C_{\Gamma_1} e^{-NM} \int d\mu(\mathbf{A}) e^{N\beta \text{Tr}[\frac{1}{N} \mathbf{T}^\top \mathbf{A}_2 \mathbf{T}]} e^{\frac{N}{2} \ln(\det(\mathbf{A}_1))} \right] \\ &= \ln C_{\Gamma_1} - NM + \ln \left[ \int d\mu(\mathbf{A}) e^{N\beta \text{Tr}[\frac{1}{N} \mathbf{T}^\top \mathbf{A}_2 \mathbf{T}]} e^{\frac{N}{2} \ln(\det(\mathbf{A}_1))} \right].\end{aligned}\quad (235)$$

3003 So long as the second term  $\text{Tr}[\mathbb{I}_M]$  does not depend on  $N$ , it will vanish when we take the partial  
3004 derivative of  $\beta\Gamma_{\mathcal{Q}^2}^{IZ}$  to obtain the  $\mathcal{E}_{gen}^{NN}$ , in which case it is not important. We can then simply write  
3005 the Generating Function  $\beta\Gamma_{\mathcal{Q}^2}^{IZ}$  as in Eqn. 118 as:

$$\beta\Gamma_{\mathcal{Q}^2}^{IZ} = \ln \left[ \int d\mu(\mathbf{A}_1) e^{N\beta \text{Tr}[\frac{1}{N} \mathbf{T}^\top \mathbf{A}_2 \mathbf{T}]} e^{\frac{N}{2} \ln(\det(\mathbf{A}_1))} \right], \quad (236)$$

3006 or, in Bra-Ket notation, as

$$\beta\Gamma_{\mathcal{Q}^2}^{IZ} = \ln \left\langle e^{N\beta \text{Tr}[\frac{1}{N} \mathbf{T}^\top \mathbf{A}_2 \mathbf{T}]} e^{\frac{N}{2} \ln(\det(\mathbf{A}_1))} \right\rangle_{\mathbf{A}}. \quad (237)$$

#### 3007 A.5 MLP3 Model Details

3008 XXX. I THINK THAT I HAVE MADE NO CHANGES TO THIS SECTION YET; PRESUMABLY  
3009 CHANGE AFTER EMPIRICAL SECTIONS ARE FINALIZED.

3010 The empirical MLP3 Model implements the assumptions described in Section 5 used the  
3011 following procedures:

3012 A three-layer Multi-Layer Perceptron was trained for classification on the MNIST dataset[?].  
3013 The first Fully Connected (FC) hidden layer has 300 units, the second FC hidden layer has 100  
3014 units, and the third FC layer has ten units for classification, matching the ten digit classes of  
3015 MNIST. Input images are grayscale, and were rescaled to the  $[0, 1]$  range. Following the keras[?]  
3016 defaults, the weights were initialized using the Glorot Normal[?] method, and the biases were  
3017 initialized to 0. Each model was trained using Categorical Cross Entropy as the loss function. The  
3018 loss function was *summed* over each mini-batch, which is the default behavior for Keras, rather  
3019 than being *averaged*, which is the default for pytorch[?].

3020 Optimization was carried out by either Stochastic Gradient Descent (SGD) without momentum,  
3021 or the Adam algorithm [?]. The Learning Rate (LR) was set to 0.01 for SGD, and 0.001 for  
3022 Adam. The LR was held constant, i.e., there was no decay schedule. Each algorithm proceeded  
3023 epoch by epoch until the value of the loss function did not decrease by more than 0.0001 for three  
3024 consecutive epochs. At each epoch, the **WeightWatcher** tool was used to compute metrics for each  
3025 layer. Loss values reported are the average loss per labeled example, and not the summed loss  
3026 over each minibatch. Training loss is averaged over all batches in the epoch, whereas test loss is  
3027 evaluated once at the end of the epoch.

3028 In some experiments, only one layer was trained, while the others were left frozen. In other  
3029 experiments all layers were trained. Models were trained using a series of mini-batch sizes ranging  
3030 from 1 to 32. For each separate training run, the models were re-initialized to the same starting  
3031 random weights, all random seeds were reset, and deterministic computations were used to train  
3032 the models.

3033 Separate notebooks are provided for keras and pytorch implementations of the experiments.

## A.6 Tanaka's Result

In this section, we will rederive the result by Tanaka [?, ?] that we use in our main derivation, and, importantly, explain how to address the missing Temperature term. For completeness, we restate it here using the notation of the main text:

$$\lim_{N \gg 1} \frac{1}{N} \ln \underbrace{\left( \exp \left( \frac{\beta}{2} \text{Tr} [\mathbf{W}^\top \mathbf{A}_2 \mathbf{W}] \right) \right)}_{\text{HCIZ Integral}} \Big|_{\mathbf{A}} = \frac{\beta}{2} \sum_{i=1}^M \mathbb{G}_{\mathbf{A}}(\lambda_i) \quad (238)$$

where  $\mathbf{W}$  is the  $N \times M$  Teacher weight matrix,  $\mathbf{A} = \mathbf{A}_2$  is the  $N \times N$  Student (correlation) matrix, but  $\beta$  is now the inverse-Temperature (because we are working with real matrices), and we have added a  $\frac{1}{2}$  (which will be clear later).  $\mathbb{G}_{\mathbf{A}}(\lambda)$  is a complex analytic function of the eigenvalues  $\lambda$  of (the Teacher Correlation matrix)  $\mathbf{X}$ , whose functional form will depend on the structure of the limiting form of (the Student) ESD  $\rho_{\mathbf{A}}^\infty(\lambda)$ . We call it the **Norm? Generating Function**; and we may also write it as  $\mathbb{G}_{\mathbf{A}}(\mathbf{X})$  below.

To apply this result, we note that while the term  $\beta$  is just a constant in [?] (1 or 2, depending on whether the random matrix is real or complex), it is not actually inverse Temperature  $\beta = \frac{1}{T}$  in the original derivation. Still, we seek a final result that is linear in  $\beta = \frac{1}{T}$ , so that we can easily evaluate  $\bar{Q}^2$  in the high-T limit, i.e.  $\bar{Q}^2 = \frac{\partial}{\partial N} \frac{1}{\beta} \beta \mathbf{\Gamma}_{\bar{Q}^2, N \gg 1}^{IZ} = \frac{\partial}{\partial \beta} \frac{1}{N} \beta \mathbf{\Gamma}_{\bar{Q}^2, N \gg 1}^{IZ}$  (see 11). We can introduce  $\beta = \frac{1}{T}$  by simply changing the scale of  $\mathbf{A}_2$  since the final result is a sum of R-transforms, which by definition are linear, i.e.  $\mathbb{G}_{\mathbf{A}}(\beta \lambda) = \beta \mathbb{G}_{\mathbf{A}}(\lambda)$ , however, it is instructive rederive the final result, with  $\beta$  explicitly included.

**Notation.** We start by re-writing the Tanaka result, Eqn. (238), in our notation for the expected value  $\langle \dots \rangle_{\mathbf{A}}$  operator, as follows:

$$\frac{1}{2} \beta \mathbf{\Gamma}_{\bar{Q}^2, N \gg 1}^{IZ} = \lim_{N \gg 1} \ln \underbrace{\int d\mu(\mathbf{A}) \left[ \exp \left( \frac{\beta}{2} \text{Tr} [\mathbf{W}^\top \mathbf{A}_2 \mathbf{W}] \right) \right]}_{\text{HCIZ Integral}} = N \beta \frac{1}{2} \sum_{i=1}^M \mathbb{G}_{\mathbf{A}}(\lambda_i). \quad (239)$$

where we have added a  $\frac{1}{2}$  for technical convenience (to make the connection with the LDP, below). If we denote the internal HCIZ integral as

$$\mathbb{Z}^{IZ} := \int d\mu(\mathbf{A}) \left[ \exp \left( \frac{\beta}{2} \text{Tr} [\mathbf{W}^\top \mathbf{A}_2 \mathbf{W}] \right) \right], \quad (240)$$

then it holds that

$$\beta \mathbf{\Gamma}_{\bar{Q}^2}^{IZ} := \ln \mathbb{Z}^{IZ}, \quad (241)$$

from which it follows that

$$\beta \mathbf{\Gamma}_{\bar{Q}^2, N \gg 1}^{IZ} := \lim_{N \gg 1} \ln \mathbb{Z}^{IZ}. \quad (242)$$

The SPA approximates the Partition Function  $\mathbb{Z}^{IZ}$ , which is now an HCIZ integral, by its peak value. For this,  $\mathbb{G}_{\mathbf{A}}(\lambda)$  itself must either not explicitly depend on  $N$  and/or at least not grow faster than  $N$ .

The trick here is we can choose an R-transform of  $\mathbf{A}$  that is a simple analytic expression based on the observed the empirical spectral density (ESD) of the  $\mathbf{X}$ . And this can readily be done for the ESDs for a wide range of layer weight matrices observed in modern DNNs because the their ESDs are Heavy-Tailed Power Law [?]. We can then readily express the Quality  $\bar{Q}$  of the Teacher layer in a simple functional form, (i.e an approximate Shatten Norm),

Importantly, the matrices  $\mathbf{X}$  and  $\mathbf{A}$  must be well approximated by low rank matrices since the derivation in Tanaka requires this. Fortunately, this appears to be generally true for the layers in very well trained DNNs, which is what allows us to apply this withing the ECS.

Finally, we note that  $\mathbb{G}_{\mathbf{A}}(\mathbf{X})$  is kind of *Generalized Norm* because it can be evaluated as a sum over a function of the  $M$  eigenvalues  $\lambda_{\mu}$  of the Teacher correlation matrix  $\mathbf{X} = \frac{1}{N} \mathbf{W}^{\top} \mathbf{W}$ .  $\mathbb{G}_{\mathbf{A}}(\mathbf{X})$  will turn out to be a simple expression similar to the Frobenius Norm or the Shatten Norm of  $\mathbf{X}$ , depending on the functional form we choose to model the limiting form of the Student ESD,  $\rho_{\mathbf{A}}^{\infty}(\lambda)$ .

### A.6.1 Setup and Outline

To evaluate 239, we want to integrate over all Student Correlation matrices  $\mathbf{A}$  that “resemble the Teacher Correlation matrix  $\mathbf{X}$ ”. To formalize this idea, we need to define the measure over “all desired  $\mathbf{A}$ ”,  $d\mu(\mathbf{A})$ , in terms of the actual  $M$  eigenvalues,  $\{\lambda_i\}_{i=1}^M$ , of the Teacher.

**Using a source matrix  $\mathbf{D}$  to represent  $d\mu(\mathbf{A})$  with  $d\mu(\mathbf{W})$ .** We consider all matrices  $\mathbf{A}$  with the same limiting spectral density,  $\rho_{\mathbf{A}}^{\infty}(\lambda)$ , as the limiting (*empirical*) ESD of the Teacher. That is, we want  $\rho_{\mathbf{A}}^{\infty}(\lambda) = \rho_{\mathbf{W}}^{\infty}(\lambda)$ , where  $\mathbf{T} = \mathbf{W}$ . **[Comment on the nature of the randomness requiring to be invariant under unitary transformations.]** Of course, there are infinitely many weight matrices  $\mathbf{W}$  with the same  $M$  eigenvalues,  $\{\lambda_i\}_{i=1}^M$ , as the Teacher. Let us specify these matrices with the measure  $d\mu(\mathbf{W})$ . Doing this lets us then write the measure  $d\mu(\mathbf{A})$  in terms of  $d\mu(\mathbf{W})$  as:

$$d\mu(\mathbf{A}) := e^{-\frac{\beta}{2} \text{Tr}[\mathbf{W} \mathbf{D} \mathbf{W}^{\top}]} d\mu(\mathbf{W}), \quad (243)$$

where  $\mathbf{D}$  is some  $M \times M$  matrix, called the Source Matrix, to be specified below, and the  $\frac{1}{2}$  here as well. Indeed, the key idea here will be to define  $\mathbf{D}$  in such a way as to obtain the desired final result. Notice also that we have added a  $\beta$  term; this will be factored out later.

We can now represent the partition function  $\mathbb{Z}^{IZ}$ , by inserting Eqn. 243 into Eqn. 240.  $\mathbb{Z}^{IZ}$  is now defined as an integral over all possible (Teacher) weight matrices  $\mathbf{W}$

$$\mathbb{Z}^{IZ} = \int d\mathbf{W} \exp\left[\frac{\beta}{2} (\text{Tr}[\mathbf{W}^{\top} \mathbf{A}_2 \mathbf{W}] - \text{Tr}[\mathbf{W} \mathbf{D} \mathbf{W}^{\top}])\right], \quad (244)$$

Observe that this integral only converges when all the eigenvalues of  $\mathbf{D}$ ,  $\{\vartheta_{\mu}\}_{\mu=1}^M$ , are larger than the maximum eigenvalue of  $\mathbf{A}$ , i.e., when  $\vartheta_{\mu} > \lambda_{\max}$ , for  $\mu \in [1, M]$  (although below this will become  $\beta \vartheta_{\mu} > \lambda_{\max}$ ). Later, we will place  $\mathbf{D}$  in diagonal form, and we will obtain an explicit expression for its  $M$  eigenvalues in terms of the  $M$  non-zero eigenvalues of  $\mathbf{X}$ . The eigenvalues of  $\mathbf{D}$  will turn out to Lagrange Multipliers, needed later.

### The Saddle Point Approximation (SPA) and the Large Deviation Principle (LDP).

To evaluate the large- $N$  case of  $\beta \Gamma_{\mathcal{Q}^2}^{IZ}$  (see 241, 242), we assume that the distribution of possible Teacher correlation matrices,  $\mu(\mathbf{X})$ , satisfies a *Large Deviation Principle (LDP)*. A LDP applies to probability distributions that take an exponential form, such that  $\mu(\mathbf{X}) = e^{-NI(\mathbf{X})} d\mu(\mathbf{X})$ , where  $I(\mathbf{X})$  is Entropy or Rate function  $I(\mathbf{X})$ . In applying a LDP, we effectively restrict measure of student correlation matrices  $\mathbf{A}$  to those most similar to the empirically observed Teacher correlation matrix  $\mathbf{X}$ . We expect the measure over all Teacher correlation matrices follows an LDP because the ESD is far from Gaussian, the dominant generalizing components reside in the tail of the ESD, and at finite-size the tail decays at worst as an exponentially Truncated Power Law (TPL).

3101 **Two steps to evaluate  $\langle \mathbb{Z}^{IZ} \rangle_{\mathbf{A}}$  in the large- $N$  approximation.** The goal is to start with  
 3102 Eqn. 244 and obtain two separate, equivalent relations, Eqns. 245 and 247:

- 3103 1. **Obtaining an integral transform of  $\rho_{\mathbf{A}}^{\infty}(\lambda)$ .** First, we expand and reduce Eqn. 244 and  
 3104 evaluate the expected value of  $\mathbb{E}_{\mathbf{A}}[\mathbb{Z}^{IZ}] = \mathbb{E}_{\mathbf{A}_2}[\mathbb{Z}^{IZ}]$  in the large- $N$  limit by expressing the  
 3105  $\rho_{\mathbf{A}}(\lambda)$  for the  $N \times N$  matrix  $\mathbf{A} = \mathbf{A}_2 = \frac{1}{N} \mathbf{S} \mathbf{S}^{\top}$  in the continuum representation, i.e., as  
 3106  $\rho_{\mathbf{A}}^{emp}(\lambda) \rightarrow \rho_{\mathbf{A}}^{\infty}(\lambda)$ , to obtain:

$$\lim_{N \gg 1} \frac{1}{N} \ln \mathbb{E}_{\mathbf{A}_2}[\mathbb{Z}^{IZ}] = M \ln\left(\frac{2\pi}{\beta}\right) - \sum_{\mu=1}^M \int \ln(\delta_{\mu} - \lambda) \rho_{\mathbf{A}}^{\infty}(\lambda) d\lambda. \quad (245)$$

3107 This gives us an  $\mathbb{E}_{\mathbf{A}_2}[\mathbb{Z}^{IZ}]$  in terms of an integral transform  $\rho_{\mathbf{A}}^{\infty}(\lambda)$ , which we can model.<sup>43</sup>

- 3108 2. **Forming the Saddle Point Approximation (SPA).** We evaluate Eqn. 244 as the expected  
 3109 value of  $\mathbb{E}_{\mathbf{A}}[\mathbb{Z}^{IZ}] = \mathbb{E}_{\mathbf{A}_1}[\mathbb{Z}^{IZ}]$  for the  $M \times M$  matrix  $\mathbf{A} = \mathbf{A}_1 = \frac{1}{N} \mathbf{S}^{\top} \mathbf{S}$  (but explicitly in terms  
 3110 of  $d\mu(\mathbf{X})$ ). Then, taking in the large- $N$  approximation using the SPA, (and which can be  
 3111 done implicitly using the LDP), we obtain

$$\lim_{N \gg 1} \mathbb{E}_{\mathbf{A}_1}[\mathbb{Z}^{IZ}] \simeq \int \exp(\beta N \text{Tr}[\mathcal{G}(\mathbf{X})]) d\mu(\mathbf{X}) \approx \exp(\beta N \mathcal{G}^{max}) \quad (246)$$

3112 where  $\mathcal{G}(\mathbf{X})$  depends on  $\mathbb{G}_{\mathbf{A}}(\mathbf{X})$ , and  $\mathcal{G}^{max} = \sup_{\mathbf{X}} \mathcal{G}(\mathbf{X})$ . We can then write

$$\lim_{N \gg 1} \frac{1}{N} \ln \mathbb{E}_{\mathbf{A}_1}[\mathbb{Z}^{IZ}] \approx \beta \mathcal{G}^{max}, \quad (247)$$

- 3113 3. **Finding the Inverse Legendre Transform.** To do this, we now equate

$$\lim_{N \gg 1} \frac{1}{N} \ln \mathbb{E}_{\mathbf{A}_1}[\mathbb{Z}^{IZ}] = \lim_{N \gg 1} \frac{1}{N} \ln \mathbb{E}_{\mathbf{A}_2}[\mathbb{Z}^{IZ}] \quad (248)$$

3114 Then, we can form the inverse Legendre transform which we will let us relate  $\mathbb{G}_{\mathbf{A}}(\lambda)$  in  
 3115 Eqn. 238 to the integrated R-transform of  $\rho_{\mathbf{A}}^{\infty}(\lambda)$ .

3116 (See A.6.5.)

### 3117 A.6.2 Step 1. Forming the Integral Transformation of ESD ( $\rho_{\mathbf{A}}^{\infty}(\lambda)$ )

3118 We first establish Eqn. 245, in Steps 1.1 – 1.4. This is done by changing variables under a Unitary  
 3119 transformation,  $\mathbf{W} \rightarrow \tilde{\mathbf{W}}$ , evaluating the resulting functional determinant, and then taking the  
 3120 continuum limit of the ESD  $\tilde{\rho}_{\mathbf{A}}(\lambda) \rightarrow \rho_{\mathbf{A}}^{\infty}(\lambda)$ .

3121 **Step 1.1** To do so, let us first assume that Teacher correlation matrix  $\mathbf{X}$  and the source matrix  
 3122  $\mathbf{D}$  are simultaneously diagonalizable (i.e., their commutator is zero:  $[\mathbf{X}, \mathbf{D}] = 0$ ). In this case, we  
 3123 may write the generating function  $\mathbb{Z}^{IZ}$  in Eqn. 244 as

$$\mathbb{Z}^{IZ} = \int d\mu(\mathbf{W}) \exp \frac{\beta}{2} \left( \text{Tr}[\mathbf{W}^{\top} \mathbf{U}^{\top} \mathbf{A} \mathbf{U} \mathbf{W}] - \text{Tr}[\mathbf{W} \mathbf{V}^{\top} \mathbf{D} \mathbf{V} \mathbf{W}^{\top}] \right), \quad (249)$$

3124 where we have defined

$$\mathbf{A}_2 = \mathbf{U}^{\top} \mathbf{A} \mathbf{U}, \quad \mathbf{D} = \mathbf{V}^{\top} \mathbf{D} \mathbf{V}, \quad (250)$$

<sup>43</sup>This integral of  $\rho_{\mathbf{A}}^{\infty}(\lambda)$  is related to the *Shannon Transform*, an integral transform from information theory that is useful when analyzing the mutual information or the capacity of a communication channel [?].



where  $\mathbf{U}$  ( $N \times N$ ) and  $\mathbf{V}$  ( $M \times M$ ) are Unitary matrices. Since  $\mathbf{U}^\top \mathbf{U} = \mathbf{I}$  and  $\mathbf{V}^\top \mathbf{V} = \mathbf{I}$ , we can insert these identities into  $\mathbb{Z}^{IZ}$  in 249, giving

$$\mathbb{Z}^{IZ} = \int d\mu(\mathbf{W}) \exp \frac{\beta}{2} \times \left( \text{Tr}[(\mathbf{V}^\top \mathbf{V}) \mathbf{W}^\top \mathbf{U}^\top \mathbf{\Lambda} \mathbf{U} \mathbf{W} (\mathbf{V}^\top \mathbf{V})] - \text{Tr}[(\mathbf{U}^\top \mathbf{U}) \mathbf{W} \mathbf{V}^\top \mathbf{\Delta} \mathbf{V} \mathbf{W}^\top (\mathbf{U}^\top \mathbf{U})] \right). \quad (251)$$

We can identify the reduced weight matrix  $\check{\mathbf{W}}$  as

$$\check{\mathbf{W}} = \mathbf{U} \mathbf{W} \mathbf{V}^\top, \quad \check{\mathbf{W}}^\top = \mathbf{V} \mathbf{W}^\top \mathbf{U}^\top, \quad (252)$$

Rearranging parentheses, this gives

$$\mathbb{Z}^{IZ} = \int d\mu(\mathbf{W}) \exp \frac{\beta}{2} \times \left( \text{Tr}[\mathbf{V}^\top (\mathbf{V} \mathbf{W}^\top \mathbf{U}^\top) \mathbf{\Lambda} (\mathbf{U} \mathbf{W} \mathbf{V}^\top) \mathbf{V}] - \text{Tr}[\mathbf{U}^\top (\mathbf{U} \mathbf{W} \mathbf{V}^\top) \mathbf{\Delta} (\mathbf{V} \mathbf{W}^\top \mathbf{U}^\top) \mathbf{U}] \right). \quad (253)$$

We can now express  $\mathbb{Z}^{IZ}$  in terms of  $\check{\mathbf{W}}$  as

$$\mathbb{Z}^{IZ} = \int d\mu(\mathbf{W}) \exp \frac{\beta}{2} \left( \text{Tr}[\mathbf{V}^\top \check{\mathbf{W}}^\top \mathbf{\Lambda} \check{\mathbf{W}} \mathbf{V}] - \text{Tr}[\mathbf{U}^\top \check{\mathbf{W}} \mathbf{\Delta} \check{\mathbf{W}}^\top \mathbf{U}] \right). \quad (254)$$

Since the Trace operator  $\text{Tr}[\cdot]$  is invariant to Unitary (Orthogonal) transformations, we can now remove the  $\mathbf{U}$  and  $\mathbf{V}$  terms, giving the simplified expression for our generating function  $\mathbb{Z}^{IZ}$  in terms of the two diagonal matrices  $\mathbf{\Lambda}, \mathbf{\Delta}$ , the reduced weight matrix  $\check{\mathbf{W}}$ , and the Jacobian  $J(\check{\mathbf{W}})$  transformation for  $d\mu(\mathbf{W}) \rightarrow d\mu(\check{\mathbf{W}})$ , as:

$$\mathbb{Z}^{IZ} = \int d\mu(\check{\mathbf{W}}) J(\check{\mathbf{W}}) \exp \frac{\beta}{2} \left( \text{Tr}[\check{\mathbf{W}}^\top \mathbf{\Lambda} \check{\mathbf{W}}] - \text{Tr}[\check{\mathbf{W}} \mathbf{\Delta} \check{\mathbf{W}}^\top] \right). \quad (255)$$

**Step 1.2** We can now evaluate the integral using the standard relation for the functional determinant for infinite-dimensional Gaussian integrals [?]

$$\mathbb{Z}^{IZ} = \left( \frac{2\pi}{\beta} \right)^{NM/2} \det(\mathbf{\Delta} - \mathbf{\Lambda})^{-1/2} \quad (256)$$

where the Jacobian is unity for the Unitary transformation.

$$J(\check{\mathbf{W}}) = 1. \quad (257)$$

since  $\mathbf{W} \mapsto \check{\mathbf{W}}$  is an orthogonal transformation. We now use the standard Trace-Log-Determinant relation [?]

$$\text{Tr}[\ln \mathbf{M}] = \ln \det \mathbf{M}. \quad (258)$$

Let us insert  $(\exp \ln)$  on the R.H.S. of 256, to obtain

$$\begin{aligned} \mathbb{Z}^{IZ} &= \exp \ln \left[ \left( \frac{2\pi}{\beta} \right)^{NM/2} \det(\mathbf{\Delta} - \mathbf{\Lambda})^{-1/2} \right] \\ &= \exp \left[ \left( \frac{NM}{2} \right) \ln \frac{2\pi}{\beta} - \frac{1}{2} \text{Tr}[\ln(\mathbf{\Delta} - \mathbf{\Lambda})] \right] \\ &= \exp \left[ \frac{NM}{2} \ln \frac{2\pi}{\beta} - \frac{1}{2} \ln \det(\mathbf{\Delta} - \mathbf{\Lambda}) \right]. \end{aligned} \quad (259)$$

3140 **Step 1.3** We now want to express the generating function  $\mathbb{Z}^{IZ}$  in 259 in terms of an integral  
 3141 over the continuous, limiting spectral density  $\rho_{\mathbf{A}}(\lambda)$  of the correlation matrix  $\mathbf{A}_2$ .

3142 First, we express the Determinant of the matrix  $\mathbf{\Delta} - \mathbf{\Lambda}$  in terms of discrete eigenvalues:

$$\det(\mathbf{\Delta} - \mathbf{\Lambda})^{-1/2} = \prod_{\mu=1}^M \prod_{i=1}^N (\vartheta_{\mu} - \lambda_i)^{-1/2}. \quad (260)$$

3143 This gives the Log-Determinant in terms of the  $M$  (non-zero) eigenvalues of  $\mathbf{D}$  and  $\mathbf{A}_2$ , as

$$\ln \det(\mathbf{\Delta} - \mathbf{\Lambda})^{-1/2} = -\frac{1}{2} \sum_{\mu=1}^M \sum_{i=1}^N \ln(\vartheta_{\mu} - \lambda_i). \quad (261)$$

3144 We can express the ESD,  $\tilde{\rho}_{\mathbf{A}}(\lambda)$ , of the Student Correlation matrix  $\mathbf{A}_2$  in terms of the Dirac  
 3145 delta-function,  $\delta(x)$ , as

$$\tilde{\rho}_{\mathbf{A}}(\lambda) = \sum_{i=1}^N \delta(\lambda - \lambda_i). \quad (262)$$

3146 [I hate that we use  $\delta$  for both the eigenvalues and as a delta function] Using this, the Expected Value  
 3147 of the Log-Determinant in 261 can be expressed in terms of the ESD of  $\mathbf{A}_2$  as

$$\begin{aligned} \left\langle \ln \det(\mathbf{\Delta} - \mathbf{\Lambda})^{-1/2} \right\rangle_{\mathbf{A}_2} &= -\frac{1}{2} \sum_{\mu=1}^M \sum_{i=1}^N \int d\lambda \ln(\vartheta_{\mu} - \lambda) \delta(\lambda - \lambda_i) \\ &= -\frac{1}{2} \sum_{\mu=1}^M \int d\lambda \ln(\vartheta_{\mu} - \lambda) \sum_{i=1}^N \delta(\lambda - \lambda_i) \\ &= -\frac{1}{2} \sum_{\mu=1}^M \int d\lambda \ln(\vartheta_{\mu} - \lambda) \tilde{\rho}_{\mathbf{A}}(\lambda). \end{aligned} \quad (263)$$

3148 [Check the normalization factor  $N$  on the outside, and the above eqn]

3149 Let us insert this back into our expression for the generating function, 259, giving  $\mathbb{E}_{\mathbf{A}_2}[\mathbb{Z}^{IZ}]$  in  
 3150 terms of the ESD  $\tilde{\rho}_{\mathbf{A}}$  as

$$\mathbb{E}_{\mathbf{A}_2}[\mathbb{Z}^{IZ}] = \exp \left\{ \frac{N}{2} \left[ M \ln \frac{2\pi}{\beta} - \sum_{\mu=1}^M \int d\lambda \ln(\vartheta_{\mu} - \lambda) \tilde{\rho}_{\mathbf{A}}(\lambda) \right] \right\}. \quad (264)$$

3151 We can now replace the sum over the  $N$  eigenvalues  $\lambda_i$  with an integral over the limiting ESD,  
 3152  $\rho(\lambda)$ , to obtain

$$\rho_{\mathbf{A}}^{\infty}(\lambda) = \lim_{N \rightarrow \infty} \tilde{\rho}_{\mathbf{A}}(\lambda). \quad (265)$$

3153 Observe that this effectively means that we are taking a large- $N$  limit,  $N \gg 1$ . This lets us write  
 3154 the Expected Value of the generating function  $\mathbb{Z}^{IZ}$  in 264 as

$$\lim_{N \gg 1} \mathbb{E}_{\mathbf{A}_2}[\mathbb{Z}^{IZ}] = \exp \left\{ \frac{N}{2} \left[ M \ln \frac{2\pi}{\beta} - \sum_{\mu=1}^M \int d\lambda \ln(\vartheta_{\mu} - \lambda) \rho_{\mathbf{A}}^{\infty}(\lambda) \right] \right\} \quad (266)$$

3155 **Step 1.4** Using the Self-Averaging Property,

$$\ln \mathbb{E}_{\mathbf{A}_2}[\mathbb{Z}^{IZ}] \simeq \langle \ln \mathbb{Z}^{IZ} \rangle_{\mathbf{A}_1}, \quad (267)$$

3156 It follows from Eqn. 266 that

$$\lim_{N \gg 1} \ln \mathbb{E}_{\mathbf{A}_2}[\mathbb{Z}^{IZ}] \simeq \frac{NM}{2} \ln \frac{2\pi}{\beta} - \frac{N}{2} \sum_{\mu=1}^M \int d\lambda \ln(\vartheta_{\mu} - \lambda) \rho_{\mathbf{A}}^{\infty}(\lambda). \quad (268)$$

3157 The  $N$ -dependence now cancels out, and we are left an approximate expression due to the  
 3158 remaining dependence of the continuum limiting density  $\rho_{\mathbf{A}}^\infty(\lambda)$  (for  $\mathbf{A} = \mathbf{A}_2$ )

$$\lim_{N \gg 1} \frac{2}{N} \ln \mathbb{E}_{\mathbf{A}_2}[\mathbb{Z}^{IZ}] = M \ln \frac{2\pi}{\beta} - \sum_{\mu=1}^M \int d\lambda \ln(\vartheta_\mu - \lambda) \rho_{\mathbf{A}}^\infty(\lambda). \quad (269)$$

3159 This completes the derivation of Eqn. 245; we have an expression for the expected value of  $\mathbb{Z}^{IZ}$ ,  
 3160 evaluated in the large- $N$  (continuum) limit.

### 3161 A.6.3 Step 2: The Saddle Point Approximation (SPA): Explicitly forming the Large 3162 Deviation Principle (LDP)

3163 We now evaluate  $\mathbb{E}_{\mathbf{A}}[\mathbb{Z}^{IZ}]$  in Eqn. 246 as  $\mathbb{E}_{\mathbf{A}_2}[\mathbb{Z}^{IZ}]$  to establish Eqn. 247, .

3164 Using the LDP (and following similar approaches in spin glass theory [?]), below we will show  
 3165 that we can write the expected value of  $\mathbb{Z}^{IZ}$  in terms of  $d\mu(\mathbf{X})$  now (which is equivalent to  $d\mu(\mathbf{A}_1)$ )  
 3166 and in the large- $N$  approximation, as

$$\lim_{N \gg 1} \mathbb{E}_{\mathbf{A}_1}[\mathbb{Z}^{IZ}] = \int \exp(\beta N \text{Tr}[\mathbb{G}_{\mathbf{A}}(\mathbf{X})] - NI(\mathbf{X}) + o(N)) d\mu(\mathbf{X}) \quad (270)$$

3167 where  $I(\mathbf{X})$  is Rate Function, defined below, and  $\mathbb{G}_{\mathbf{A}}(\mathbf{X})$  is what we are eventually solving for.

3168 We first introduce a new change of measure,  $d\mu(\mathbf{W}) \rightarrow d\mu(\mathbf{X})$ . Then, we show this lets us  
 3169 express  $\mathbb{E}_{\mathbf{A}_1}[\mathbb{Z}^{IZ}]$  as  $\mathbb{E}_{\mathbf{X}}[\mathbb{Z}^{IZ}]$  and to express it using the LDP. Next, we apply a SPA to solve for  
 3170  $\mathcal{G}^{max}$ . Importantly, we also show how to incorporate the inverse-Temperature  $\beta$ , which is new.

3171

3172 **Step 2.1** To define the transformation  $d\mu(\mathbf{W}) \rightarrow d\mu(\mathbf{X})$ , where (recall)  $\mathbf{X} = \frac{1}{N} \mathbf{W}^\top \mathbf{W}$ , we use  
 3173 the (again) the integral representation of the Dirac delta-function  $\delta(x)$ :

$$\delta(x) := \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{i\hat{x}x} d\hat{x}. \quad (271)$$

3174 This lets us express the transformation of measure  $d\mu(\mathbf{W}) \rightarrow d\mu(\mathbf{X})$  (approximately) as

$$\begin{aligned} d\mu(\mathbf{W}) &:= \delta\left(\frac{1}{2} \text{Tr}[N\mathbf{X} - \mathbf{W}^\top \mathbf{W}]\right) d\mu(\mathbf{X}) \\ &= \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{i\frac{1}{2} \text{Tr}[\hat{X}(N\mathbf{X} - \mathbf{W}^\top \mathbf{W})]} d\mu(\mathbf{X}) d\hat{X}, \end{aligned} \quad (272)$$

3175 where  $\hat{X}$  is a scalar (or really matrix of scalars), and we have a 1/2 term for mathematical  
 3176 consistency below. <sup>44</sup>

3177 **Step 2.2** Next, we take a Wick Rotation,  $i\hat{\mathbf{X}} \rightarrow \hat{\mathbf{X}}$ , so that the terms under the integral are all  
 3178 real (not complex), giving:

$$d\mu(\mathbf{W}) = \frac{1}{2\pi} \int_{-i\infty}^{i\infty} e^{\frac{1}{2} \text{Tr}[\hat{X}(N\mathbf{X} - \mathbf{W}^\top \mathbf{W})]} d\mu(\mathbf{X}) d\hat{X}. \quad (273)$$

3179 [Did we flip a sign here. Do we have  $i\hat{\mathbf{X}} \rightarrow -\hat{\mathbf{X}}$  ?]

<sup>44</sup>The full change of measure would require a delta function constraint for each matrix element  $X_{i,j}$ , i.e.,  $\delta(\frac{1}{2}N(X_{i,j} - [\mathbf{W}^\top \mathbf{W}]_{i,j}))$ . Here, we assume the Trace constraint is sufficient for our level of rigor.

3180 **Step 2.3** We now insert 273 into 244, which lets express  $\mathbb{Z}^{IZ}$  as an integral over the Teacher  
 3181 Correlation matrices

$$\begin{aligned}\mathbb{Z}^{IZ} &= \frac{1}{2\pi} \int_{\mathbf{X}} \int_{-\infty}^{i\infty} e^{N\frac{\beta}{2} \text{Tr}[\mathcal{G}(\mathbf{X})] + \text{Tr}[\hat{X}\mathbf{X}]} e^{-\frac{1}{2} \text{Tr}[\hat{X}\mathbf{W}^\top \mathbf{W}]} e^{\frac{\beta}{2} \text{Tr}[\mathbf{W}\mathbf{D}\mathbf{W}^\top]} d\hat{X} d\mu(\mathbf{X}) \\ &= \frac{1}{2\pi} \int_{\mathbf{X}} \int_{-\infty}^{i\infty} e^{N\frac{\beta}{2} \text{Tr}[\mathcal{G}(\mathbf{X})] + \frac{1}{2} \text{Tr}[\hat{X}\mathbf{X}]} e^{-\frac{1}{2} \text{Tr}[\mathbf{W}\hat{X}\mathbf{W}^\top] + \frac{\beta}{2} \text{Tr}[\mathbf{W}\mathbf{D}\mathbf{W}^\top]} d\hat{X} d\mu(\mathbf{X}) \\ &= \frac{1}{2\pi} \int_{\mathbf{X}} \int_{-\infty}^{i\infty} e^{N\frac{\beta}{2} \text{Tr}[\mathcal{G}(\mathbf{X})] + \frac{1}{2} \text{Tr}[\hat{X}\mathbf{X}]} e^{\frac{1}{2} \text{Tr}[\mathbf{W}(\beta\mathbf{D} - \hat{X})\mathbf{W}^\top]} d\hat{X} d\mu(\mathbf{X}).\end{aligned}\quad (274)$$

3182 [ABOVE may be missing a 1/2]

3183 **Step 2.4** We can now rearrange terms to make this expression look like the Eqn. 270  
 3184 In Large Deviations Theory, the Rate Function is defined by the Legendre Transform,

$$\mathcal{I}(\mathbf{X}) = \sup_{\tilde{\mathbf{X}}} \left[ \text{Tr} \frac{1}{2} \mathbf{X}^\top \tilde{\mathbf{X}} - \ln \mathbb{M}(\tilde{\mathbf{X}}) \right], \quad (275)$$

3185 where  $\mathbb{M}(\tilde{\mathbf{X}})$  is the Moment Generating Function,  $\ln \mathbb{M}(\tilde{\mathbf{X}})$ , is the Cumulant Generating Function,  
 3186 and  $\tilde{\mathbf{X}}$  is a (matrix of) *Lagrange Multiplier(s)*.  $\mathbb{M}(\tilde{\mathbf{X}})$  is defined in terms of the (unnormalized)  
 3187 density  $p(\mathbf{x})$  as

$$\mathbb{M}(\tilde{\mathbf{X}}) = \int \exp \left( \frac{1}{2} \mathbf{x}^\top \tilde{\mathbf{X}} \mathbf{x} \right) p(\mathbf{x}) d\mathbf{x} \quad (276)$$

3188 which, in term, is defined in terms of the source matrix  $\mathbf{D}$ ,

$$p(\mathbf{x}) = \exp \left( -\frac{1}{2} \mathbf{x}^\top \beta \mathbf{D} \mathbf{x} \right). \quad (277)$$

3189 The moment generating function  $\mathbb{M}(\tilde{\mathbf{X}})$  is then given by

$$\mathbb{M}(\tilde{\mathbf{X}}) = \int \exp \left( -\frac{1}{2} \mathbf{x}^\top (\beta \mathbf{D} - \tilde{\mathbf{X}}) \mathbf{x} \right) d\mathbf{x} = (2\pi)^{\frac{M}{2}} \det(\beta \mathbf{D} - \tilde{\mathbf{X}})^{-\frac{1}{2}}. \quad (278)$$

3190 **Step 2.5** The Saddle Point Approximation (SPA) can be used to solve for  $\mathcal{I}(\tilde{\mathbf{X}})$  by solving for  
 3191 the stationary conditions

$$\frac{\partial}{\partial \tilde{\mathbf{X}}} \mathcal{I}(\mathbf{X}, \tilde{\mathbf{X}}) = 0. \quad (279)$$

3192 First, let us compute  $\ln \mathbb{M}(\tilde{\mathbf{X}})$  as:

$$\ln \mathbb{M}(\tilde{\mathbf{X}}) = \frac{M}{2} \ln(2\pi) - \frac{1}{2} \ln \det(\beta \mathbf{D} - \tilde{\mathbf{X}}). \quad (280)$$

3193 Substituting this into the expression for the Legendre transform, we obtain:

$$\mathcal{I}(\mathbf{X}, \tilde{\mathbf{X}}) = \sup_{\tilde{\mathbf{X}}} \left[ \frac{1}{2} \text{Tr}[\mathbf{X}\tilde{\mathbf{X}}] - \frac{M}{2} \ln(2\pi) + \frac{1}{2} \ln \det(\beta \mathbf{D} - \tilde{\mathbf{X}}) \right]. \quad (281)$$

3194 The supremum of this expression is attained at the value of  $\tilde{\mathbf{X}}$  that satisfies:

$$\frac{\partial}{\partial \tilde{\mathbf{X}}} \left[ \frac{1}{2} \text{Tr}[\mathbf{X}\tilde{\mathbf{X}}] + \frac{1}{2} \ln \det(\beta \mathbf{D} - \tilde{\mathbf{X}}) \right] = 0. \quad (282)$$

3195 Taking the derivative, we obtain

$$\frac{1}{2} \mathbf{X} + \frac{1}{2} (\beta \mathbf{D} - \tilde{\mathbf{X}})^{-1} = 0, \quad (283)$$

3196 which simplifies to:

$$\mathbf{X} = (\beta \mathbf{D} - \tilde{\mathbf{X}})^{-1} \Rightarrow \tilde{\mathbf{X}} = \beta \mathbf{D} - \mathbf{X}^{-1}. \quad (284)$$

3197 Substituting  $\tilde{\mathbf{X}} = \beta \mathbf{D} - \mathbf{X}^{-1}$  back into the expression for  $I(\mathbf{X})$ , we obtain:

$$I(\mathbf{X}) = \frac{1}{2} \left[ \text{Tr} [\mathbf{X}(\beta \mathbf{D} - \mathbf{X}^{-1})] - \frac{M}{2} \ln(2\pi) + \frac{1}{2} \ln \det (\mathbf{X}^{-1}) \right]. \quad (285)$$

3198

$$\text{Tr} [\mathbf{X}\beta \mathbf{D} - \mathbf{I}] = \text{Tr} [\mathbf{X}\beta \mathbf{D}] - N, \quad (286)$$

3199

$$\ln \det (\mathbf{X}^{-1}) = -\ln \det (\mathbf{X}), \quad (287)$$

3200 we get:

$$I(\mathbf{X}) = \frac{1}{2} [\text{Tr} [\mathbf{X}\beta \mathbf{D}] - \ln \det (\mathbf{X}) - M - M \ln(2\pi)]. \quad (288)$$

3201 Finally, we express  $I(\mathbf{X})$  in the form:

$$I(\mathbf{X}) = \frac{1}{2} [-M(1 + \ln(2\pi)) + \text{Tr} [\mathbf{X}\beta \mathbf{D}] - \ln \det (\mathbf{X})]. \quad (289)$$

3202 **[THE DERIVATION ABOVE FOR  $\mathcal{G}(\mathbf{X})$  may have some TYPOS: CHECK]**

## Step 2.6

$$\beta \mathcal{G}(\mathbf{X}) = M(1 + \ln 2\pi) + \beta \text{Tr} [\mathbb{G}_{\mathbf{A}}(\mathbf{X})] - \text{Tr} [\mathbf{X}\beta \mathbf{D}] + \ln \det (\mathbf{X}). \quad (290)$$

3203 We restrict our solution to those where  $\mathbf{X}$  and  $\beta \mathbf{D}$  can be diagonalized simultaneously. In  
3204 particular, this lets us write

$$\text{Tr} [\mathbf{X}\beta \mathbf{D}] = \sum_{\mu=1}^M \beta \delta_{\mu} \lambda_{\mu}, \quad (291)$$

3205 where  $\beta \delta_{\mu}$  and  $\lambda_{\mu}$  denote the eigenvalues of  $\mathbf{X}$  and  $\beta \mathbf{D}$ , resp.

3206 We can now write the maximum value of  $\mathbb{G}_{\mathbf{A}}$ ,  $\mathcal{G}^{max}$ , as

$$\beta \mathcal{G}^{max} = M \left( 1 + \ln \frac{2\pi}{\beta} \right) - \sum_{\mu=1}^M \min_{\beta \delta_{\mu}} [\beta \delta_{\mu} \lambda_{\mu} - \beta \mathbb{G}_{\mathbf{A}}(\lambda_{\mu}) + \ln \lambda_{\mu}]. \quad (292)$$

## 3207 A.6.4 Expressing the Norm Generating Function ( $\mathbb{G}_{\mathbf{A}}(\lambda)$ ) as the Integrated R- 3208 transform ( $R(z)$ ) of the Correlation Matrix ( $\mathbf{A}$ )

3209 Having completed both steps, let us combine Eqns. 245, 269 with 247 and 292. We follow the first  
3210 arguments by Tanaka[?] (which follows Cherrier[?]).

$$M \ln \left( \frac{2\pi}{\beta} \right) - \sum_{\mu=1}^M \int \ln(\beta \delta_{\mu} - \lambda) \rho_{\mathbf{A}}^{\infty}(\lambda) d\lambda = M \left( 1 + \ln \frac{2\pi}{\beta} \right) - \sum_{\mu=1}^M \min_{\beta \delta_{\mu}} [\beta \delta_{\mu} \lambda_{\mu} - \beta \mathbb{G}_{\mathbf{A}}(\lambda_{\mu}) + \ln \lambda_{\mu}]. \quad (293)$$

3211 By cancelling the  $\ln \frac{2\pi}{\beta}$  term from both sides, we obtain

$$- \sum_{\mu=1}^M \int \ln(\beta \delta_{\mu} - \lambda) \rho_{\mathbf{A}}^{\infty}(\lambda) d\lambda = M - \sum_{\mu=1}^M \min_{\beta \delta_{\mu}} [\beta \delta_{\mu} \lambda_{\mu} - \beta \mathbb{G}_{\mathbf{A}}(\lambda_{\mu}) + \ln \lambda_{\mu}]. \quad (294)$$

3212 Since this is true for every  $\mu$ , we can solve this for any arbitrary eigenvalue  $\lambda_{\mu}$ .

3213 Dropping the  $\mu$  subscript, we have the following identity:

$$\min_{\delta} [\beta\delta\lambda - \beta\mathbb{G}_{\mathbf{A}}(\lambda) + \ln \lambda] = 1 - \int \ln(\beta\delta - \lambda) \rho_{\mathbf{A}}^{\infty}(\lambda) d\lambda. \quad (295)$$

3214 We need to invert 295 in order to find  $\beta\mathbb{G}_{\mathbf{A}}(\lambda)$ . If we choose the eigenvalues of  $\mathbf{D}$  such that  
 3215  $\beta\delta_{\mu} > \lambda_{max}$  for all  $\mu$ , then this relation is concave and therefore invertible via a Legendre transform.

3216 This gives

$$\beta\mathbb{G}_{\mathbf{A}}(\lambda) = \beta\delta(\lambda)\lambda - \int \ln[\beta\delta(\lambda) - \lambda] \rho_{\mathbf{A}}^{\infty}(\lambda) d\lambda - \ln \lambda - 1, \quad (296)$$

3217 where we need to define  $\beta\delta(\lambda)$ , which (not to be confused with the Dirac delta-function), describes  
 3218 the functional dependence between the eigenvalues of the source matrix  $\mathbf{D}$  and the Student  
 3219 Correlation Matrix  $\mathbf{A}$ .

3220  $\mathbb{G}_{\mathbf{A}}(\lambda)$  is computed by minimizing over  $\delta$ , ensuring the relationship holds for the entire spectrum.  
 3221 So let us take the derivative of  $\beta\mathbb{G}_{\mathbf{A}}$  w/r.t  $\lambda$ . Term by term, this gives:

$$\frac{d}{d\lambda} \beta\delta(\lambda)\lambda = \beta\delta(\lambda) + \frac{d\beta\delta(\lambda)}{d\lambda} \lambda \quad (297)$$

3222

$$\frac{d}{d\lambda} \ln \lambda = \frac{1}{\lambda} \quad (298)$$

3223

$$\begin{aligned} \frac{d}{d\lambda} \int \ln[\beta\delta(\lambda) - \lambda] \rho_{\mathbf{A}}^{\infty}(\lambda) d\lambda &= \int \frac{d}{d\lambda} \ln[\beta\delta(\lambda) - \lambda] \rho_{\mathbf{A}}^{\infty}(\lambda) d\lambda = \int \frac{d\beta\delta(\lambda)}{d\lambda} \frac{\rho_{\mathbf{A}}^{\infty}(\lambda)}{\beta\delta(\lambda) - \lambda} d\lambda \\ &= \frac{d\beta\delta(\lambda)}{d\lambda} \int \frac{\rho_{\mathbf{A}}^{\infty}(\lambda)}{\beta\delta(\lambda) - \lambda} d\lambda \end{aligned} \quad (299)$$

3224 We can now simplify by defining  $\delta(\lambda)$  implicitly by the integral relation

$$\lambda = \int \frac{\rho_{\mathbf{A}}^{\infty}(\lambda)}{\beta\delta(\lambda) - \lambda} d\lambda. \quad (300)$$

3225 Combining terms, this gives

$$\frac{d\beta\mathbb{G}_{\mathbf{A}}(\lambda)}{d\lambda} = \beta\delta(\lambda) - \frac{1}{\lambda}, \quad (301)$$

3226 Inverting the derivative, we obtain an integral equation for  $\beta\mathbb{G}_{\mathbf{A}}(\lambda)$

$$\beta\mathbb{G}_{\mathbf{A}}(\lambda) = \int_0^{\lambda} \left( \beta\delta(z) - \frac{1}{z} \right) dz. \quad (302)$$

3227 Notice since  $\beta\delta(\lambda) \approx \frac{1}{\lambda}$  for  $\lambda \ll 1$ , then as  $\mathbb{G}_{\mathbf{A}}(0) = 0$  and we set the lower integrand to 0 (for  
 3228 now).

3229 **THIS LAST PART MAY NEED SOME DISCUSSION: SEE TAnAKA 2007**

3230 To further connect these results with the Cauchy transform  $\mathcal{C}_{\mathbf{A}}(z)$  and the R-transform  $R_{\mathbf{A}}(z)$ ,  
 3231 we recall that the Cauchy transform is given by: **[this is also defined in the main text; need to make**  
 3232 **sure we have the same sign convention]**

$$\mathcal{C}_{\mathbf{A}}(z) = \int \frac{\rho_{\mathbf{A}}(\lambda)}{z - \lambda} d\lambda. \quad (303)$$

3233 The relationship between the Cauchy transform and the R-transform is expressed as:

$$\mathcal{C}_{\mathbf{A}} \left( R_{\mathbf{A}}(z) + \frac{1}{z} \right) = z, \quad (304)$$

3234 which implies:

$$\beta \mathbb{G}_{\mathbf{A}}(\lambda) = \int_0^\lambda R_{\mathbf{A}}(z) dz. \quad (305)$$

3235 Although we note that, at least for our purposes,  $R(z)$  may and probably will have a branchcut  
3236 at the start of the tail of ESD of  $\rho_{\mathbf{A}}$ , so we actually want

$$\beta \mathbb{G}_{\mathbf{A}}(\lambda) = \int_{\lambda_{min}^{ECS}}^\lambda R_{\mathbf{A}}(z) dz. \quad (306)$$

3237 where  $\lambda_{min}^{ECS}$  corresponds to the start of the branchcut if necessary.

### 3238 **A.6.5 Selecting $\mathbf{A} := \mathbf{A}_1$ instead of $\mathbf{A}_2$**

3239 In principle, we could have selected  $\mathbf{A} := \mathbf{A}_1 = \frac{1}{N} \mathbf{S}^\top \mathbf{S}$  for the Student Correlation matrix, thereby  
3240 avoiding the discussion on the Duality of Measures altogether. Doing this, however, would make  
3241  $\mathbf{A}$   $M \times M$ , thereby require defining the Source Matrix  $\mathbf{D}$  as an  $N \times N$  matrix, with presumably  
3242  $N - M$  zero eigenvalues. This would cause  $\mathbf{D}$  to violate the condition  $\vartheta_\mu > \beta\lambda$  for all eigenvalues  $\lambda$   
3243 of  $\mathbf{A}_1$ . In this case, it would be challenging to define the large- $N$  limit.

## 3244 **A.7 The Inverse-Wishart (IW) Model**

3245 In this section, we rederive the integral  $G(\lambda)[IW]$  for the Inverse Wishart (IW) model, focusing  
3246 on the branch cut starting at  $z = \kappa/2$  and extending to infinity. This branch cut corresponds to  
3247 the support of the ESD in this region. We will:

- 3248 1. Explain the presence of the branch cut and its implications.
- 3249 2. Show that  $R(z)[IW]$  becomes complex along this branch cut because the term under the  
3250 square root becomes negative.
- 3251 3. Perform the integral  $G(\lambda)[IW]$ , showing all steps.
- 3252 4. Compute the modulus  $|G(\lambda)[IW]| = \sqrt{G(\lambda)[IW]^* G(\lambda)[IW]}$  to obtain a real-valued esti-  
3253 mate, analogous to a probability estimate.

### 3254 **A.7.1 The Branch Cut in the IW Model**

3255 The R-transform for the IW model is given by:

$$R(z)[IW] = \frac{\kappa - \sqrt{\kappa(\kappa - 2z)}}{z}, \quad (307)$$

3256 where  $\kappa > 0$  is a parameter related to the dimensions of the random matrices under consideration.  
3257 The function  $\sqrt{\kappa(\kappa - 2z)}$  introduces a branch point at  $z = \kappa/2$  because the argument of the square  
3258 root becomes zero at this point:

$$\kappa - 2z = 0 \quad \Rightarrow \quad z = \frac{\kappa}{2}. \quad (308)$$

3259 For  $z > \kappa/2$ , the argument  $\kappa - 2z$  becomes negative, and thus the square root becomes imaginary.  
3260 This leads to a branch cut starting at  $z = \kappa/2$  and extending to  $z = \infty$  along the real axis. This  
3261 branch cut affects the analyticity of  $R(z)[IW]$ , and it must be carefully considered in the integral  
3262  $G(\lambda)[IW]$ .



### 3263 **A.7.2** $R(z)[IW]$ is Complex Along the Branch Cut

3264 For  $z > \kappa/2$ , we have:

$$\kappa - 2z < 0 \Rightarrow \sqrt{\kappa(\kappa - 2z)} = \sqrt{-\kappa(2z - \kappa)} = i\sqrt{\kappa(2z - \kappa)}. \quad (309)$$

3265 Therefore,  $R(z)[IW]$  becomes complex:

$$R(z)[IW] = \frac{\kappa - i\sqrt{\kappa(2z - \kappa)}}{z} = \frac{\kappa}{z} - i\frac{\sqrt{\kappa(2z - \kappa)}}{z}. \quad (310)$$

3266 This expression shows that  $R(z)[IW]$  has both real and imaginary parts when  $z > \kappa/2$ .

### 3267 **A.7.3** Calculation of $G(\lambda)[IW]$

3268 We aim to compute the integral:

$$G(\lambda)[IW] = \int_{z_0}^{\lambda} R(z)[IW] dz, \quad (311)$$

3269 where  $z_0 \geq \kappa/2$ .

3270 **Integrating the Real Part.** First, we consider the real part of  $R(z)[IW]$ :

$$\operatorname{Re}[R(z)[IW]] = \frac{\kappa}{z}. \quad (312)$$

3271 The integral of this real part is:

$$G_{\text{real}}(\lambda)[IW] = \int_{z_0}^{\lambda} \frac{\kappa}{z} dz = \kappa [\ln z]_{z_0}^{\lambda} = \kappa (\ln \lambda - \ln z_0). \quad (313)$$

3272 **Integrating the Imaginary Part.** Next, consider the imaginary part:

$$\operatorname{Im}[R(z)[IW]] = -\frac{\sqrt{\kappa(2z - \kappa)}}{z}. \quad (314)$$

3273 If we let  $u = 2z - \kappa$ , then:

$$z = \frac{u + \kappa}{2}, \quad dz = \frac{du}{2}. \quad (315)$$

3274 Substituting this into the imaginary part:

$$\operatorname{Im}[R(z)[IW]] = -\frac{\sqrt{\kappa u}}{\frac{u + \kappa}{2}} = -\frac{2\sqrt{\kappa u}}{u + \kappa}, \quad (316)$$

3275 the integral becomes:

$$G_{\text{imag}}(\lambda)[IW] = \int_{u_0}^{u_{\lambda}} -\frac{2\sqrt{\kappa u}}{u + \kappa} \cdot \frac{du}{2} = -\int_{u_0}^{u_{\lambda}} \frac{\sqrt{\kappa u}}{u + \kappa} du, \quad (317)$$

3276 where  $u_0 = 2z_0 - \kappa$  and  $u_{\lambda} = 2\lambda - \kappa$ . If we simplify the integrand:

$$\sqrt{\kappa u} = \sqrt{\kappa} \sqrt{u}, \quad (318)$$

3277 then the integral becomes:

$$G_{\text{imag}}(\lambda)[IW] = -\sqrt{\kappa} \int_{u_0}^{u_{\lambda}} \frac{\sqrt{u}}{u + \kappa} du. \quad (319)$$

3278 This integral can be evaluated using standard integral formulas. We will compute it step by step.

3279 **Evaluating the Integral.** Consider the integral:

$$I = \int \frac{\sqrt{u}}{u + \kappa} du. \quad (320)$$

3280 We can use the following integral formula:

$$\int \frac{\sqrt{u}}{u + a} du = 2\sqrt{u} - 2a \tan^{-1} \left( \frac{\sqrt{u}}{\sqrt{a}} \right) + C, \quad (321)$$

3281 where  $a > 0$  and  $u > 0$ . Applying this formula, we get:

$$I = 2\sqrt{u} - 2\kappa \tan^{-1} \left( \frac{\sqrt{u}}{\sqrt{\kappa}} \right) + C. \quad (322)$$

3282 Therefore, the imaginary part of  $G(\lambda)[IW]$  is:

$$\begin{aligned} G_{\text{imag}}(\lambda)[IW] &= -\sqrt{\kappa} \left[ 2\sqrt{u} - 2\kappa \tan^{-1} \left( \frac{\sqrt{u}}{\sqrt{\kappa}} \right) \right] u_0^{u\lambda} \\ &= -\sqrt{\kappa} \left( \left[ 2\sqrt{u_\lambda} - 2\kappa \tan^{-1} \left( \frac{\sqrt{u_\lambda}}{\sqrt{\kappa}} \right) \right] - \left[ 2\sqrt{u_0} - 2\kappa \tan^{-1} \left( \frac{\sqrt{u_0}}{\sqrt{\kappa}} \right) \right] \right) \\ &= -2\sqrt{\kappa} (\sqrt{u_\lambda} - \sqrt{u_0}) + 2\kappa^{3/2} \left( \tan^{-1} \left( \frac{\sqrt{u_\lambda}}{\sqrt{\kappa}} \right) - \tan^{-1} \left( \frac{\sqrt{u_0}}{\sqrt{\kappa}} \right) \right). \end{aligned} \quad (323)$$

3283 **Combining Real and Imaginary Parts.** Combine the real and imaginary parts to obtain  
3284  $G(\lambda)[IW]$ :

$$G(\lambda)[IW] = G_{\text{real}}(\lambda)[IW] + iG_{\text{imag}}(\lambda)[IW]. \quad (324)$$

3285 Substituting the expressions:

$$\begin{aligned} G(\lambda)[IW] &= \kappa (\ln \lambda - \ln z_0) \\ &\quad + i \left( -2\sqrt{\kappa} (\sqrt{u_\lambda} - \sqrt{u_0}) + 2\kappa^{3/2} \left( \tan^{-1} \left( \frac{\sqrt{u_\lambda}}{\sqrt{\kappa}} \right) - \tan^{-1} \left( \frac{\sqrt{u_0}}{\sqrt{\kappa}} \right) \right) \right). \end{aligned} \quad (325)$$

3286 Recall that  $u = 2z - \kappa$ , so:

$$\sqrt{u} = \sqrt{2z - \kappa}. \quad (326)$$

3287 Therefore, we can write  $G(\lambda)[IW]$  as:

$$\begin{aligned} G(\lambda)[IW] &= \kappa \ln \left( \frac{\lambda}{z_0} \right) - 2i\sqrt{\kappa} (\sqrt{2\lambda - \kappa} - \sqrt{2z_0 - \kappa}) \\ &\quad + 2i\kappa^{3/2} \left( \tan^{-1} \left( \frac{\sqrt{2\lambda - \kappa}}{\sqrt{\kappa}} \right) - \tan^{-1} \left( \frac{\sqrt{2z_0 - \kappa}}{\sqrt{\kappa}} \right) \right). \end{aligned} \quad (327)$$

#### 3288 **A.7.4 Computing the Modulus $|G(\lambda)[IW]|$**

3289 To obtain a real-valued estimate, we compute the modulus of  $G(\lambda)[IW]$ :

$$|G(\lambda)[IW]| = \sqrt{(\text{Re}[G(\lambda)[IW]])^2 + (\text{Im}[G(\lambda)[IW]])^2}. \quad (328)$$

3290 **Calculating the Real Part Square.** The real part is:

$$\text{Re}[G(\lambda)[IW]] = \kappa \ln \left( \frac{\lambda}{z_0} \right). \quad (329)$$

3291 Therefore,

$$(\text{Re}[G(\lambda)[IW]])^2 = \kappa^2 \left( \ln \left( \frac{\lambda}{z_0} \right) \right)^2. \quad (330)$$

3292 **Calculating the Imaginary Part Square.** The imaginary part is:

$$\text{Im}[G(\lambda)[IW]] = -2\sqrt{\kappa} \left( \sqrt{2\lambda - \kappa} - \sqrt{2z_0 - \kappa} \right) + 2\kappa^{3/2} \left( \tan^{-1} \left( \frac{\sqrt{2\lambda - \kappa}}{\sqrt{\kappa}} \right) - \tan^{-1} \left( \frac{\sqrt{2z_0 - \kappa}}{\sqrt{\kappa}} \right) \right). \quad (331)$$

3293 Let's denote:

$$A = -2\sqrt{\kappa} \left( \sqrt{2\lambda - \kappa} - \sqrt{2z_0 - \kappa} \right), \quad B = 2\kappa^{3/2} \left( \tan^{-1} \left( \frac{\sqrt{2\lambda - \kappa}}{\sqrt{\kappa}} \right) - \tan^{-1} \left( \frac{\sqrt{2z_0 - \kappa}}{\sqrt{\kappa}} \right) \right). \quad (332)$$

3294 Then,

$$(\text{Im}[G(\lambda)[IW]])^2 = (A + B)^2 = A^2 + 2AB + B^2. \quad (333)$$

3295 **Computing the Modulus.** The modulus is:

$$|G(\lambda)[IW]| = \sqrt{\left( \kappa \ln \left( \frac{\lambda}{z_0} \right) \right)^2 + (A + B)^2}. \quad (334)$$

3296 **Interpretation.** While the expression for  $|G(\lambda)[IW]|$  appears complex, it encapsulates the  
 3297 cumulative effect of both the real and imaginary components of  $G(\lambda)[IW]$ . This modulus provides  
 3298 a real-valued estimate that is meaningful in the context of probability estimates.

### 3299 **A.7.5 Summary**

3300 By integrating  $R(z)[IW]$  directly, including its complex components, we have obtained an explicit  
 3301 expression for  $G(\lambda)[IW]$  as a complex function. Computing the modulus  $|G(\lambda)[IW]|$  gives us a  
 3302 real-valued function that accurately captures the contribution of the tail of the ESD in the IW  
 3303 model. This approach accounts for the complex nature of  $R(z)[IW]$  along the branch cut  $z > \kappa/2$ ,  
 3304 and it provides a meaningful point estimate for further analysis.