

## The Principles of Deep Learning Theory

This textbook establishes a theoretical framework for understanding deep learning models of practical relevance. With an approach that borrows from theoretical physics, Roberts and Yaida provide clear and pedagogical explanations of how realistic deep neural networks actually work. To make results from the theoretical forefront accessible, the authors eschew the subject's traditional emphasis on intimidating formality without sacrificing accuracy. Straightforward and approachable, this volume balances detailed first-principle derivations of novel results with insight and intuition for theorists and practitioners alike. This self-contained textbook is ideal for students and researchers interested in artificial intelligence with minimal prerequisites of linear algebra, calculus, and informal probability theory, and it can easily fill a semester-long course on deep learning theory. For the first time, the exciting practical advances in modern artificial intelligence capabilities can be matched with a set of effective principles, providing a timeless blueprint for theoretical research in deep learning.

**Daniel A. Roberts** was cofounder and CTO of Diffeo, an AI company acquired by Salesforce; a research scientist at Facebook AI Research; and a member of the School of Natural Sciences at the Institute for Advanced Study in Princeton, NJ. He was a Hertz Fellow, earning a PhD from MIT in theoretical physics, and was also a Marshall Scholar at Cambridge and Oxford Universities.

**Sho Yaida** is a research scientist at Meta AI. Prior to joining Meta AI, he obtained his PhD in physics at Stanford University and held postdoctoral positions at MIT and at Duke University. At Meta AI, he uses tools from theoretical physics to understand neural networks, the topic of this book.

**Boris Hanin** is an assistant professor at Princeton University in the Operations Research and Financial Engineering Department. Prior to joining Princeton in 2020, Boris was an assistant professor at Texas A&M in the Math Department and an NSF postdoc at MIT. He has taught graduate courses on the theory and practice of deep learning at both Texas A&M and Princeton.

### Prepublication praise

“In the history of science and technology, the engineering artifact often comes first: the telescope, the steam engine, digital communication. The theory that explains its function and its limitations often appears later: the laws of refraction, thermodynamics, and information theory. With the emergence of deep learning, AI-powered engineering wonders have entered our lives — but our theoretical understanding of the power and limits of deep learning is still partial. This is one of the first books devoted to the theory of deep learning, and lays out the methods and results from recent theoretical approaches in a coherent manner.”

– **Prof. Yann LeCun**, *New York University and Chief AI Scientist at Meta*

“For a physicist, it is very interesting to see deep learning approached from the point of view of statistical physics. This book provides a fascinating perspective on a topic of increasing importance in the modern world.”

– **Prof. Edward Witten**, *Institute for Advanced Study*

“This is an important book that contributes big, unexpected new ideas for unraveling the mystery of deep learning’s effectiveness, in unusually clear prose. I hope it will be read and debated by experts in all the relevant disciplines.”

– **Prof. Scott Aaronson**, *University of Texas at Austin*

“It is not an exaggeration to say that the world is being revolutionized by deep learning methods for AI. But why do these deep networks work? This book offers an approach to this problem through the sophisticated tools of statistical physics and the renormalization group. The authors provide an elegant guided tour of these methods, interesting for experts and non-experts alike. They write with clarity and even moments of humor. Their results, many presented here for the first time, are the first steps in what promises to be a rich research program, combining theoretical depth with practical consequences.”

– **Prof. William Bialek**, *Princeton University*

“This book’s physics-trained authors have made a cool discovery, that feature learning depends critically on the ratio of depth to width in the neural net.”

– **Prof. Gilbert Strang**, *Massachusetts Institute of Technology*

# The Principles of Deep Learning Theory

An Effective Theory Approach  
to Understanding Neural Networks

DANIEL A. ROBERTS  
*MIT*

SHO YAIDA  
*Meta AI*

*based on research in collaboration with*

BORIS HANIN  
*Princeton University*



CAMBRIDGE  
UNIVERSITY PRESS

University Printing House, Cambridge CB2 8BS, United Kingdom  
One Liberty Plaza, 20th Floor, New York, NY 10006, USA  
477 Williamstown Road, Port Melbourne, VIC 3207, Australia  
314–321, 3rd Floor, Plot 3, Splendor Forum, Jasola District Centre, New Delhi – 110025, India  
103 Penang Road, #05–06/07, Visioncrest Commercial, Singapore 238467

Cambridge University Press is part of the University of Cambridge.

It furthers the University's mission by disseminating knowledge in the pursuit of education, learning, and research at the highest international levels of excellence.

[www.cambridge.org](http://www.cambridge.org)

Information on this title: [www.cambridge.org/9781316519332](http://www.cambridge.org/9781316519332)

DOI: 10.1017/9781009023405

© Cambridge University Press 2022

This publication is in copyright. Subject to statutory exception and to the provisions of relevant collective licensing agreements, no reproduction of any part may take place without the written permission of Cambridge University Press.

First published 2022

*A catalogue record for this publication is available from the British Library.*

*Library of Congress Cataloging-in-Publication Data*

Names: Roberts, Daniel A., 1987– author.

Title: The principles of deep learning theory : an effective theory approach to understanding neural networks / Daniel A. Roberts and Sho Yaida based on research in collaboration with Boris Hanin.

Description: New York : Cambridge University Press, 2022. | Includes bibliographical references and index.

Identifiers: LCCN 2021060635 (print) | LCCN 2021060636 (ebook) |

ISBN 9781316519332 (hardback) | ISBN 9781009023405 (epub)

Subjects: LCSH: Deep learning (Machine learning) |

BISAC: SCIENCE / Physics / Mathematical & Computational

Classification: LCC Q325.73 .R63 2022 (print) | LCC Q325.73 (ebook) |

DDC 006.3/1–dc23/eng20220215

LC record available at <https://lcn.loc.gov/2021060635>

LC ebook record available at <https://lcn.loc.gov/2021060636>

ISBN 9781316519332 Hardback

Additional resources for this publication at [www.cambridge.org/deeplearningtheory](http://www.cambridge.org/deeplearningtheory)

Cambridge University Press has no responsibility for the persistence or accuracy of URLs for external or third-party internet websites referred to in this publication and does not guarantee that any content on such websites is, or will remain, accurate or appropriate.

# Contents

<b>Preface</b>	<b>ix</b>
<b>0 Initialization</b>	<b>1</b>
0.1 An Effective Theory Approach . . . . .	2
0.2 The Theoretical Minimum . . . . .	3
<b>1 Pretraining</b>	<b>11</b>
1.1 Gaussian Integrals . . . . .	12
1.2 Probability, Correlation and Statistics, and All That . . . . .	21
1.3 Nearly-Gaussian Distributions . . . . .	26
<b>2 Neural Networks</b>	<b>37</b>
2.1 Function Approximation . . . . .	37
2.2 Activation Functions . . . . .	43
2.3 Ensembles . . . . .	47
<b>3 Effective Theory of Deep Linear Networks at Initialization</b>	<b>53</b>
3.1 Deep Linear Networks . . . . .	54
3.2 Criticality . . . . .	56
3.3 Fluctuations . . . . .	59
3.4 Chaos . . . . .	65
<b>4 RG Flow of Preactivations</b>	<b>71</b>
4.1 First Layer: Good-Old Gaussian . . . . .	73
4.2 Second Layer: Genesis of Non-Gaussianity . . . . .	79
4.3 Deeper Layers: Accumulation of Non-Gaussianity . . . . .	90
4.4 Marginalization Rules . . . . .	96
4.5 Subleading Corrections . . . . .	100
4.6 RG Flow and RG Flow . . . . .	103
<b>5 Effective Theory of Preactivations at Initialization</b>	<b>109</b>
5.1 Criticality Analysis of the Kernel . . . . .	110
5.2 Criticality for Scale-Invariant Activations . . . . .	123
5.3 Universality Beyond Scale-Invariant Activations . . . . .	125

5.3.1	General Strategy . . . . .	126
5.3.2	No Criticality: Sigmoid, Softplus, Nonlinear Monomials, etc. . . . .	128
5.3.3	$K^* = 0$ Universality Class: tanh, sin, etc. . . . .	130
5.3.4	Half-Stable Universality Classes: SWISH, etc. and GELU, etc. . . . .	135
5.4	Fluctuations . . . . .	137
5.4.1	Fluctuations for the Scale-Invariant Universality Class . . . . .	139
5.4.2	Fluctuations for the $K^* = 0$ Universality Class . . . . .	141
5.5	Finite-Angle Analysis for the Scale-Invariant Universality Class . . . . .	146
<b>6</b>	<b>Bayesian Learning</b>	<b>153</b>
6.1	Bayesian Probability . . . . .	154
6.2	Bayesian Inference and Neural Networks . . . . .	156
6.2.1	Bayesian Model Fitting . . . . .	157
6.2.2	Bayesian Model Comparison . . . . .	165
6.3	Bayesian Inference at Infinite Width . . . . .	169
6.3.1	The Evidence for Criticality . . . . .	169
6.3.2	Let's Not Wire Together . . . . .	173
6.3.3	Absence of Representation Learning . . . . .	178
6.4	Bayesian Inference at Finite Width . . . . .	179
6.4.1	Hebbian Learning, Inc. . . . .	179
6.4.2	Let's Wire Together . . . . .	182
6.4.3	Presence of Representation Learning . . . . .	186
<b>7</b>	<b>Gradient-Based Learning</b>	<b>191</b>
7.1	Supervised Learning . . . . .	192
7.2	Gradient Descent and Function Approximation . . . . .	194
<b>8</b>	<b>RG Flow of the Neural Tangent Kernel</b>	<b>199</b>
8.0	Forward Equation for the NTK . . . . .	200
8.1	First Layer: Deterministic NTK . . . . .	206
8.2	Second Layer: Fluctuating NTK . . . . .	207
8.3	Deeper Layers: Accumulation of NTK Fluctuations . . . . .	211
8.3.0	<i>Interlude: Interlayer Correlations</i> . . . . .	211
8.3.1	NTK Mean . . . . .	215
8.3.2	NTK–Preactivation Cross Correlations . . . . .	216
8.3.3	NTK Variance . . . . .	221
<b>9</b>	<b>Effective Theory of the NTK at Initialization</b>	<b>227</b>
9.1	Criticality Analysis of the NTK . . . . .	228
9.2	Scale-Invariant Universality Class . . . . .	233
9.3	$K^* = 0$ Universality Class . . . . .	236
9.4	Criticality, Exploding and Vanishing Problems, and None of That . . . . .	241

<b>10 Kernel Learning</b>	<b>247</b>
10.1 A Small Step . . . . .	248
10.1.1 No Wiring . . . . .	250
10.1.2 No Representation Learning . . . . .	250
10.2 A Giant Leap . . . . .	252
10.2.1 Newton's Method . . . . .	253
10.2.2 Algorithm Independence . . . . .	257
10.2.3 <i>Aside</i> : Cross-Entropy Loss . . . . .	259
10.2.4 Kernel Prediction . . . . .	261
10.3 Generalization . . . . .	264
10.3.1 Bias–Variance Tradeoff and Criticality . . . . .	267
10.3.2 Interpolation and Extrapolation . . . . .	277
10.4 Linear Models and Kernel Methods . . . . .	282
10.4.1 Linear Models . . . . .	282
10.4.2 Kernel Methods . . . . .	284
10.4.3 Infinite-Width Networks as Linear Models . . . . .	287
<b>11 Representation Learning</b>	<b>291</b>
11.1 Differential of the Neural Tangent Kernel . . . . .	293
11.2 RG Flow of the dNTK . . . . .	296
11.2.0 Forward Equation for the dNTK . . . . .	297
11.2.1 First Layer: Zero dNTK . . . . .	299
11.2.2 Second Layer: Nonzero dNTK . . . . .	300
11.2.3 Deeper Layers: Growing dNTK . . . . .	301
11.3 Effective Theory of the dNTK at Initialization . . . . .	310
11.3.1 Scale-Invariant Universality Class . . . . .	312
11.3.2 $K^* = 0$ Universality Class . . . . .	314
11.4 Nonlinear Models and Nearly-Kernel Methods . . . . .	317
11.4.1 Nonlinear Models . . . . .	318
11.4.2 Nearly-Kernel Methods . . . . .	324
11.4.3 Finite-Width Networks as Nonlinear Models . . . . .	330
<b><math>\infty</math> The End of Training</b>	<b>335</b>
$\infty.1$ Two More Differentials . . . . .	337
$\infty.2$ Training at Finite Width . . . . .	347
$\infty.2.1$ A Small Step Following a Giant Leap . . . . .	351
$\infty.2.2$ Many Many Steps of Gradient Descent . . . . .	358
$\infty.2.3$ Prediction at Finite Width . . . . .	373
$\infty.3$ RG Flow of the ddNTKs: The Full Expressions . . . . .	384
<b><math>\varepsilon</math> Epilogue: Model Complexity from the Macroscopic Perspective</b>	<b>389</b>

<b>A Information in Deep Learning</b>	<b>399</b>
A.1 Entropy and Mutual Information . . . . .	400
A.2 Information at Infinite Width: Criticality . . . . .	409
A.3 Information at Finite Width: Optimal Aspect Ratio . . . . .	411
<b>B Residual Learning</b>	<b>425</b>
B.1 Residual Multilayer Perceptrons . . . . .	428
B.2 Residual Infinite Width: Criticality Analysis . . . . .	429
B.3 Residual Finite Width: Optimal Aspect Ratio . . . . .	431
B.4 Residual Building Blocks . . . . .	436
<b>References</b>	<b>439</b>
<b>Index</b>	<b>445</b>



# Preface

*This has necessitated a complete break from the historical line of development, but this break is an advantage through enabling the approach to the new ideas to be made as direct as possible.*

P. A. M. Dirac in the 1930 preface of *The Principles of Quantum Mechanics* [1].

This is a research monograph in the style of a textbook about the theory of deep learning. While this book might look a little different from the other deep learning books that you've seen before, we assure you that it is appropriate for everyone with knowledge of linear algebra, multivariable calculus, and informal probability theory, and with a healthy interest in neural networks. Practitioner and theorist alike, we want all of you to enjoy this book. Now, let us tell you some things.

First and foremost, in this book we've strived for pedagogy in every choice we've made, placing intuition above formality. This doesn't mean that calculations are incomplete or sloppy; quite the opposite, we've tried to provide full details of every calculation – of which there are certainly very many – and place a particular emphasis on the tools needed to carry out related calculations of interest. In fact, understanding how the calculations are done is as important as knowing their results, and thus often our pedagogical focus is on the details therein.

Second, while we present the details of all our calculations, we've kept the experimental confirmations to the privacy of our own computerized notebooks. Our reason for this is simple: while there's much to learn from explaining a derivation, there's not much more to learn from printing a verification plot that shows two curves lying on top of each other. Given the simplicity of modern deep-learning packages and the availability of compute, it's easy to verify any formula on your own; we certainly have thoroughly checked them all this way, so if knowledge of the existence of such plots is comforting to you, know at least that they do exist on our personal and cloud-based hard drives.

Third, our main focus is on realistic models that are used by the deep learning community in practice: we want to study *deep* neural networks. In particular, this means that (i) a number of special results on single-hidden-layer networks will not be discussed and (ii) the *infinite-width limit* of a neural network – which is equivalent to a zero-hidden-layer network – will be introduced only as a starting point. All such idealized models will eventually be *perturbed* until they correspond to a real model. We certainly acknowledge that there's a vibrant community of deep-learning theorists devoted to

exploring different kinds of idealized theoretical limits. However, our interests are fixed firmly on providing explanations for the tools and approaches used by practitioners, in an effort to shed light on what makes them work so well.

Fourth, a large part of the book is focused on deep multilayer perceptrons. We made this choice in order to pedagogically illustrate the power of the effective theory framework – not due to any technical obstruction – and along the way we give pointers for how this formalism can be extended to other architectures of interest. In fact, we expect that many of our results have a broad applicability, and we’ve tried to focus on aspects that we expect to have lasting and universal value to the deep learning community.

Fifth, while much of the material is novel and appears for the first time in this book, and while much of our framing, notation, language, and emphasis breaks with the historical line of development, we’re also very much indebted to the deep learning community. With that in mind, throughout the book we will try to reference important prior contributions, with an emphasis on recent seminal deep-learning results rather than on being completely comprehensive. Additional references for those interested can easily be found within the work that we cite.

Sixth, this book initially grew out of a research project in collaboration with Boris Hanin. To account for his effort and then support, we’ve accordingly commemorated him on the cover. More broadly, we’ve variously appreciated the artwork, discussions, encouragement, epigraphs, feedback, management, refereeing, reintroduction, and support from Rafael Araujo, Léon Bottou, Paul Dirac, Ethan Dyer, John Frank, Ross Girshick, Vince Higgs, Yoni Kahn, Yann LeCun, Kyle Mahowald, Eric Mintun, Xiaoliang Qi, Mike Rabbat, David Schwab, Stephen Shenker, Eva Silverstein, PJ Steiner, DJ Strouse, and Jesse Thaler. Organizationally, we’re grateful to FAIR and Facebook, Diffeo and Salesforce, MIT and IAIFI, and Cambridge University Press and the arXiv.

Seventh, given intense (and variously uncertain) spacetime and energy-momentum commitment that writing this book entailed, Dan is grateful to Aya, Lumi, and Lisa Yaida; from the dual sample-space perspective, Sho is grateful to Adrienne Rothschilds and would be retroactively grateful to any hypothetical future Mark or Emily that would have otherwise been thanked in this paragraph.

Eighth, we hope that this book spreads our optimism that it *is* possible to have a general theory of deep learning, one that’s both derived from first principles and at the same time focused on describing how realistic models actually work: nearly-simple phenomena in practice should correspond to nearly-simple effective theories. We dream that this type of thinking will not only lead to more *[redacted]* AI models but also guide us toward a unifying framework for understanding universal aspects of intelligence.

As if that eightfold way of prefacing the book wasn’t nearly-enough already, please note: this book has a website, [deeplearningtheory.com](http://deeplearningtheory.com), and you may want to visit it in order to determine whether the error that you just discovered is already common knowledge. If it’s not, please let us know.

*Dan Roberts & Sho Yaida  
Remotely Located*