

NEURAL CONNECTIVITY LEARNING FOR VISUAL RECOGNITION

A Thesis

Submitted to the Faculty

in partial fulfillment of the requirements for the

degree of

Doctor of Philosophy

in

Computer Science

by

Karim Ahmed

Guarini School of Graduate and Advanced Studies

DARTMOUTH COLLEGE

Hanover, New Hampshire

August 2019

Examining Committee:

(Chair) Lorenzo Torresani

Andrew Campbell

Rogério Schmidt Feris

Bo Zhu

F. Jon Kull, Ph.D.

Dean of the Guarini School of Graduate and Advanced Studies

Abstract

Deep neural networks have emerged as one of the most prominent models for problems that require the learning of complex functions and that involve large amounts of training data such as visual recognition problems. The power of deep learning stems from the ability to learn representations optimized for a specific task, as opposed to relying on hand-crafted features. To yield good results, however, deep models often require manual architecture search which is still challenging, time-consuming and error-prone process even for expert human designers. The difficulty lies in the various design choices such as depth, filter sizes, number of feature maps, and connectivity patterns that affect the performance of the system and often require tedious periods of trial and error to identify a satisfactory set of choices for a specific task and dataset. Thus, one may argue that deep learning has replaced hand-crafted features with hand-crafted models equipped by representation learning. Consequently, there has been an increasing effort to develop Neural Architecture Search methods for automatic model selection. Unfortunately, most of the Neural Architecture Search methods are computationally expensive due to the combinatorial explosion of design options in the search space. In this thesis, we propose Neural Connectivity Learning

as a middle-ground approach between two extremes: hand-crafting the network versus learning the full architecture. The proposed Neural Connectivity Learning methods aim at bounding the search space and allowing fast architecture search, by learning the connectivity between preset network components. Specifically, in the domain of large-scale image classification, we introduce two different approaches that learn separate features of coarse image classes by learning the connectivity between fine-grained classes in convolutional neural networks. Subsequently, we present an algorithm that learns the connections between the modules of convolutional neural networks, instead of being chosen a priori by the human designer. Finally, to alleviate the computational complexity of capsule neural networks, we present a new mechanism for learning the connectivity and routing between capsules. The empirical studies conducted on several image classification datasets show that the proposed connectivity learning approaches outperform the baseline approaches that rely on fixed connectivity rules.

Contents

Abstract	ii
1 INTRODUCTION	1
1.1 Neural Architecture Search (NAS)	3
1.2 Neural Connectivity Learning (NCL)	6
1.3 Main Contributions and Outline of the Thesis	7
2 RELATED WORK	12
2.1 Learning connectivity between image classes	12
2.2 Learning connectivity between layers	16
2.3 Learning connectivity between capsules	18
3 NETWORK OF EXPERTS FOR LARGE-SCALE IMAGE CATEGORIZATION	21
3.1 Overview	21
3.2 Motivation	22
3.3 Technical approach	26
3.3.1 Learning the Generalist	27
3.3.2 Training the Network of Experts	31
3.4 Experiments	32
3.4.1 Model Analysis on CIFAR-100	32
3.4.2 Categorization on ImageNet	43
3.5 Supplementary Information	45

3.6	Summary	61
4	IMAGE CATEGORIZATION WITH LEARNED EXPERT BRANCH CONNECTIONS	62
4.1	Overview	62
4.2	Motivation	63
4.3	Technical Approach	66
4.3.1	The architecture of BRANCHCONNECT	66
4.3.2	Training BRANCHCONNECT	68
4.3.3	Inference with BRANCHCONNECT	71
4.4	Experiments	72
4.4.1	Reshaping a CNN with BRANCHCONNECT	72
4.4.2	Evaluation on CIFAR-100	73
4.4.3	Evaluation on CIFAR-10	80
4.4.4	Evaluation on ImageNet	81
4.4.5	Evaluation on Synth	81
4.5	Supplementary Information	83
4.6	Summary	88
5	CONNECTIVITY LEARNING IN MODULAR NETWORKS	97
5.1	Overview	97
5.2	Motivation	98
5.3	Technical Approach	101
5.3.1	Modular architecture	101
5.3.2	Masked architecture	103
5.3.3	MASKCONNECT: learning to connect	104
5.3.4	MASKCONNECT applied to ResNet	107
5.3.5	MASKCONNECT applied to multi-branch ResNeXt	107

5.4	Experiments	109
5.4.1	Evaluation on CIFAR-100	111
5.4.2	Evaluation on ImageNet	116
5.5	Supplementary Information	117
5.6	Summary	124
6	LEARNING THE CONNECTIVITY IN CAPSULE NETWORKS	126
6.1	Overview	126
6.2	Motivation	127
6.3	Background	130
6.3.1	Capsule Neural Networks	130
6.3.2	Attention vs. Dynamic Routing	132
6.4	STARCAPS Architecture	133
6.4.1	Overview	135
6.4.2	Attention Estimator	137
6.4.3	Straight-Through Router	137
6.5	Experiments	140
6.5.1	Experimental Setup	140
6.5.2	Evaluation on MNIST	141
6.5.3	Evaluation on SmallNorb	142
6.5.4	Evaluation on CIFAR-10/CIFAR-100	144
6.5.5	Evaluation on ImageNet	144
6.6	Summary	145
7	CONCLUSIONS AND DISCUSSION	146

List of Tables

Table 3.1	NoF _E : Evaluation on CIFAR-100	35
Table 3.2	NoF _E : Comparison of different specialties on CIFAR-100	37
Table 3.3	NoF _E : Comparison of different models on CIFAR-100	37
Table 3.4	NoF _E : Depth and model capacity vs. specialization	38
Table 3.5	NoF _E : Evaluation on ImageNet	46
Table 3.6	NoF _E : Learned specialties by the generalist from CIFAR-100 . .	50
Table 3.7	NoF _E : More learned specialties by the generalist from CIFAR-100	50
Table 3.8	NoF _E : Learned specialties by the generalist from ImageNet . .	50
Table 3.9	NoF _E : AlexNet-C100 (trained on CIFAR-100)	55
Table 3.10	NoF _E : AlexNet-Quick-C100 (trained on CIFAR-100)	56
Table 3.11	NoF _E : VGG11-C100 (trained on CIFAR-100)	57
Table 3.12	NoF _E : NIN-C100 (trained on CIFAR-100)	58
Table 3.13	NoF _E : ResNet56-C100 (trained on CIFAR-100)	59
Table 3.14	NoF _E : AlexNet-Caffe (trained on ImageNet)	60
Table 4.1	BRANCHCONNECT: Evaluation on CIFAR-100	74
Table 4.2	BRANCHCONNECT: Evaluation on CIFAR-10	80
Table 4.3	BRANCHCONNECT: Evaluation on ImageNet	82
Table 4.4	BRANCHCONNECT: Evaluation on Synth dataset	82
Table 4.5	BRANCHCONNECT: Classification accuracy on CIFAR-100	86
Table 4.6	BRANCHCONNECT: CIFAR-100 AlexNet-Quick	88

Table 4.7	BRANCHCONNECT: CIFAR-100 AlexNet-Full	89
Table 4.8	BRANCHCONNECT: CIFAR-100 Network In Network	90
Table 4.9	BRANCHCONNECT: CIFAR-100 ResNet56-4X	91
Table 4.10	BRANCHCONNECT: CIFAR-100 ResNet56	92
Table 4.11	BRANCHCONNECT: ImageNet AlexNet	93
Table 4.12	BRANCHCONNECT: ImageNet ResNet-50	94
Table 4.13	BRANCHCONNECT: ImageNet ResNet-101	95
Table 4.14	BRANCHCONNECT: DICT+2-90k	96
Table 5.1	MASKCONNECT: Evaluation on CIFAR-100 based on ResNet . . .	112
Table 5.2	MASKCONNECT: Evaluation on CIFAR-100 based on ResNeXt . .	116
Table 5.3	MASKCONNECT: Evaluation on ImageNet based on ResNeXt . .	117
Table 5.4	MASKCONNECT: CIFAR-10 accuracies based on ResNeXt	118
Table 5.5	MASKCONNECT: Mini-ImageNet accuracies based on ResNeXt .	119
Table 5.6	MASKCONNECT: CIFAR-10/100 architectures based on ResNeXt .	121
Table 5.7	MASKCONNECT: Mini-ImageNet architectures	125
Table 6.1	STARCAPS: Sensitivity to the predefined number of capsules . .	142
Table 6.2	STARCAPS: Detection of novel viewpoints on SmallNorb	143

List of Figures

Figure 1.1	Hand-crafted features vs. Learned features	2
Figure 1.2	Automated Machine Learning methods	3
Figure 1.3	Examples of different architecture search spaces	6
Figure 3.1	NoE: Network of Experts	24
Figure 3.2	NoE: Generalist accuracy on ImageNet	46
Figure 3.3	NoE: Nearest Neighbor retrieval on CIFAR-100	48
Figure 3.4	NoE: Nearest Neighbor search on ImageNet	49
Figure 3.5	NoE: Distribution of specialty sizes learned by elasso	52
Figure 4.1	BRANCHCONNECT: Overview of the architecture	66
Figure 4.2	BRANCHCONNECT: Effect of different number of branches	77
Figure 4.3	BRANCHCONNECT: Train loss vs. Test loss on CIFAR-100	80
Figure 4.4	BRANCHCONNECT: Increasing the network depth	81
Figure 5.1	MASKCONNECT: Application to modular networks	102
Figure 5.2	MASKCONNECT: Learned active connections based on ResNet	111
Figure 5.3	MASKCONNECT: Learned active connections based on ResNeXt	113
Figure 5.4	MASKCONNECT: Fixed connectivity vs. Learned connectivity	114
Figure 5.5	MASKCONNECT: Active branches vs. module depth (CIFAR-100)	120
Figure 5.6	MASKCONNECT: Active branches vs. module depth (ImageNet)	120
Figure 6.1	STARCAPS: Traditional Neural Layers vs. Capsule Layers	131

Figure 6.2	STARCAPS: Dynamic routing vs. Attention-based routing . . .	132
Figure 6.3	STARCAPS: Convolutional capsule layer in STARCAPS	134
Figure 6.4	STARCAPS: STARCAPS vs. EMCaps models trained on MNIST .	141

Bibliography

- [1] Ryan Prescott Adams, Hanna M. Wallach, and Zoubin Ghahramani. Learning the structure of deep sparse graphical models. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics, AISTATS 2010, Chia Laguna Resort, Sardinia, Italy, May 13-15, 2010*, pages 1–8, 2010. 17
- [2] Karim Ahmed, Mohammad Haris Baig, and Lorenzo Torresani. Network of experts for large-scale image categorization. In *European Conference on Computer Vision (ECCV)*, 2016. 21, 73, 74, 75, 82, 86
- [3] Karim Ahmed and Lorenzo Torresani. Branchconnect: Large-scale visual recognition with learned branch connections. *arXiv preprint arXiv:1704.06010*, 2017. 62
- [4] Karim Ahmed and Lorenzo Torresani. Connectivity learning in multi-branch networks. *arXiv preprint arXiv:1709.09582*, 2017. 97
- [5] Karim Ahmed and Lorenzo Torresani. Branchconnect: Image categorization with learned branch connections. In *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1244–1253. IEEE, 2018. 62
- [6] Karim Ahmed and Lorenzo Torresani. Maskconnect: Connectivity learning by gradient descent. In *European Conference on Computer Vision (ECCV)*, 2018. 97

- [7] Karim Ahmed and Lorenzo Torresani. Star-caps: Capsule network with straight-through attentive routing. In *Advances in Neural Information Processing Systems*, page (to appear), 2019. 126
- [8] Rahaf Aljundi, Punarjay Chakravarty, and Tinne Tuytelaars. Expert gate: Lifelong learning with a network of experts. *CoRR*, abs/1611.06194, 2016. 15
- [9] Amjad Almahairi, Nicolas Ballas, Tim Cooijmans, Yin Zheng, Hugo Larochelle, and Aaron Courville. Dynamic capacity networks. In *International Conference on Machine Learning*, pages 2549–2558, 2016. 18
- [10] Peter J Angeline, Gregory M Saunders, and Jordan B Pollack. An evolutionary algorithm that constructs recurrent neural networks. *IEEE transactions on Neural Networks*, 5(1):54–65, 1994. 4
- [11] Bowen Baker, Otkrist Gupta, Nikhil Naik, and Ramesh Raskar. Designing neural network architectures using reinforcement learning. *arXiv preprint arXiv:1611.02167*, 2016. 4
- [12] Gabriel Bender, Pieter-Jan Kindermans, Barret Zoph, Vijay Vasudevan, and Quoc Le. Understanding and simplifying one-shot architecture search. In *International Conference on Machine Learning*, pages 549–558, 2018. 4, 5
- [13] Rodrigo Benenson. Classification datasets results. http://rodrigob.github.io/are_we_there_yet/build/classification_datasets_results.html. 38, 41
- [14] Emmanuel Bengio, Pierre-Luc Bacon, Joelle Pineau, and Doina Precup. Conditional computation in neural networks for faster models. *arXiv preprint arXiv:1511.06297*, 2015. 18
- [15] Samy Bengio, Jason Weston, and David Grangier. Label embedding trees for large multi-class tasks. In *Advances in Neural Information Processing Systems 23: 24th Annual Conference on Neural Information Processing Systems 2010. Proceedings of a meeting held 6-9 December 2010, Vancouver, British Columbia, Canada.*, pages 163–171, 2010. 13

-
- [16] Yoshua Bengio. Deep learning of representations: Looking forward. *CoRR*, abs/1305.0445, 2013. 15
- [17] Yoshua Bengio. Deep learning of representations: Looking forward. In *International Conference on Statistical Language and Speech Processing*, pages 1–37. Springer, 2013. 18
- [18] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013. 18, 20, 128, 139
- [19] Alessandro Bergamo and Lorenzo Torresani. Meta-class features for large-scale object categorization on a budget. In *2012 IEEE Conference on Computer Vision and Pattern Recognition, Providence, RI, USA, June 16-21, 2012*, pages 3085–3092, 2012. 29
- [20] James Bergstra and Yoshua Bengio. Random search for hyper-parameter optimization. *Journal of Machine Learning Research*, 13:281–305, 2012. 16
- [21] James S Bergstra, Rémi Bardenet, Yoshua Bengio, and Balázs Kégl. Algorithms for hyper-parameter optimization. In *Advances in neural information processing systems*, pages 2546–2554, 2011. 4
- [22] I. Biederman. Recognition-by-components: a theory of human understanding. *Psychological Review*, 94:115–147, 1987. 22
- [23] Andrew Brock, Theodore Lim, James M Ritchie, and Nick Weston. Smash: one-shot model architecture search through hypernetworks. *arXiv preprint arXiv:1708.05344*, 2017. 4
- [24] Han Cai, Jiacheng Yang, Weinan Zhang, Song Han, and Yong Yu. Path-level network transformation for efficient architecture search. *arXiv preprint arXiv:1806.02639*, 2018. 4
- [25] Xiaojun Chang, Feiping Nie, Zhigang Ma, and Yi Yang. Balanced k-means and min-cut clustering. *CoRR*, abs/1411.6235, 2014. 29, 30

- [26] Yunpeng Chen, Yannis Kalantidis, Jianshu Li, Shuicheng Yan, and Jiashi Feng. Multi-fiber networks for video recognition. In *European Conference on Computer Vision (ECCV)*, 2018. 99
- [27] KyungHyun Cho and Yoshua Bengio. Exponentially increasing the capacity-to-computation ratio for conditional computation in deep learning. *arXiv preprint arXiv:1406.7362*, 2014. 18
- [28] François Chollet. Xception: Deep learning with depthwise separable convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1251–1258, 2017. 19, 137
- [29] Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. Fast and accurate deep network learning by exponential linear units (elus). *CoRR*, abs/1511.07289, 2015. 38, 42
- [30] Matthieu Courbariaux, Yoshua Bengio, and Jean-Pierre David. Binaryconnect: Training deep neural networks with binary weights during propagations. In *Advances in Neural Information Processing Systems 28, Montreal, Quebec, Canada*, pages 3123–3131, 2015. 69, 70, 71, 105, 106
- [31] Andrew Davis and Itamar Arel. Low-rank approximations for conditional feedforward computation in deep neural networks. *arXiv preprint arXiv:1312.4461*, 2013. 18
- [32] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Fei-Fei Li. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2009)*, 20–25 June 2009, Miami, Florida, USA, pages 248–255, 2009. 18, 32, 43, 63, 72, 81, 82, 109, 117, 140, 144
- [33] Jia Deng, Jonathan Krause, Alexander C. Berg, and Fei-Fei Li. Hedging your bets: Optimizing accuracy-specificity trade-offs in large scale visual recognition. In 2012

- IEEE Conference on Computer Vision and Pattern Recognition, Providence, RI, USA, June 16-21, 2012*, pages 3450–3457, 2012. 13
- [34] Ludovic Denoyer and Patrick Gallinari. Deep sequential neural network. *arXiv preprint arXiv:1410.0510*, 2014. 18
- [35] David Eigen, Marc’Aurelio Ranzato, and Ilya Sutskever. Learning factored representations in a deep mixture of experts. *arXiv preprint arXiv:1312.4314*, 2013. 18
- [36] Thomas Elsken, Jan-Hendrik Metzen, and Frank Hutter. Simple and efficient architecture search for convolutional neural networks. *arXiv preprint arXiv:1711.04528*, 2017. 4
- [37] Thomas Elsken, Jan Hendrik Metzen, and Frank Hutter. Neural architecture search: A survey. *arXiv preprint arXiv:1808.05377*, 2018. 3, 4, 5
- [38] Thomas Elsken, Jan Hendrik Metzen, and Frank Hutter. Efficient multi-objective neural architecture search via lamarckian evolution. *ICLR*, 2019. 4
- [39] Dumitru Erhan, Yoshua Bengio, Aaron C. Courville, Pierre-Antoine Manzagol, Pascal Vincent, and Samy Bengio. Why does unsupervised pre-training help deep learning? *Journal of Machine Learning Research*, 11:625–660, 2010. 12
- [40] Dumitru Erhan, Yoshua Bengio, Aaron C. Courville, Pierre-Antoine Manzagol, Pascal Vincent, and Samy Bengio. Why does unsupervised pre-training help deep learning? *Journal of Machine Learning Research*, 11:625–660, 2010. 79
- [41] Chrisantha Fernando, Dylan Banarse, Charles Blundell, Yori Zwols, David Ha, Andrei A. Rusu, Alexander Pritzel, and Daan Wierstra. Pathnet: Evolution channels gradient descent in super neural networks. *CoRR*, abs/1701.08734, 2017. 16
- [42] Dario Floreano, Peter Dürri, and Claudio Mattiussi. Neuroevolution: from architectures to learning. *Evolutionary Intelligence*, 1(1):47–62, 2008. 16

-
- [43] Tianshi Gao and Daphne Koller. Discriminative learning of relaxed hierarchy for large-scale visual recognition. In *ICCV*, 2011. 13
- [44] Xavier Gastaldi. Shake-shake regularization. *CoRR*, abs/1705.07485, 2017. 17
- [45] Xavier Glorot, Antoine Bordes, and Yoshua Bengio. Deep sparse rectifier neural networks. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics, AISTATS 2011, Fort Lauderdale, USA, April 11-13, 2011*, pages 315–323, 2011. 12, 57, 99
- [46] Gregory Griffin and Pietro Perona. Learning and using taxonomies for fast visual categorization. In *2008 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2008), 24-26 June 2008, Anchorage, Alaska, USA, 2008*. 13
- [47] E. J. Gumbel. Statistical theory of extreme values and some practical applications: a series of lectures. In *US Govt. Print. Office*, number 33, 1954. 139
- [48] Yiwen Guo, Anbang Yao, and Yurong Chen. Dynamic network surgery for efficient dnns. In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pages 1379–1387, 2016. 17
- [49] Yunhui Guo, Honghui Shi, Abhishek Kumar, Kristen Grauman, Tajana Rosing, and Rogerio Feris. Spottune: transfer learning through adaptive fine-tuning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4805–4814, 2019. 20
- [50] Song Han, Huizi Mao, and William J. Dally. Deep compression: Compressing deep neural network with pruning, trained quantization and huffman coding. In *International Conference on Learning Representations (ICLR)*, 2015. 17
- [51] Song Han, Jeff Pool, Sharan Narang, Huizi Mao, Shijian Tang, Erich Elsen, Bryan Catanzaro, John Tran, and William J. Dally. DSD: regularizing deep neural networks

- with dense-sparse-dense training flow. In *International Conference on Learning Representations (ICLR)*, 2016. 17
- [52] Song Han, Jeff Pool, John Tran, and William J. Dally. Learning both weights and connections for efficient neural network. In *Advances in Neural Information Processing Systems 28, Montreal, Quebec, Canada*, pages 1135–1143, 2015. 17
- [53] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017. 127
- [54] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Spatial pyramid pooling in deep convolutional networks for visual recognition. In *Computer Vision - ECCV 2014 - 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part III*, pages 346–361, 2014. 12
- [55] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *CoRR*, abs/1512.03385, 2015. 26, 41, 54, 102, 112
- [56] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *2015 IEEE International Conference on Computer Vision, ICCV 2015, Santiago, Chile, December 7-13, 2015*, pages 1026–1034, 2015. 12, 59, 63, 91, 92, 94, 95
- [57] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Computer Vision and Pattern Recognition (CVPR), 2016 IEEE Conference on*, 2016. 1, 4, 15, 63, 73, 74, 75, 80, 82, 85, 86, 99, 103, 107, 109, 119, 122, 127, 144, 148
- [58] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In *Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part IV*, pages 630–645, 2016. 101

- [59] Geoffrey E Hinton, Sara Sabour, and Nicholas Frosst. Matrix capsules with em routing. *ICLR*, 2018. 10, 18, 19, 127, 128, 129, 130, 131, 132, 133, 135, 137, 140, 141, 142, 143, 144, 149
- [60] Geoffrey E. Hinton, Nitish Srivastava, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Improving neural networks by preventing co-adaptation of feature detectors. *CoRR*, abs/1207.0580, 2012. 12
- [61] Geoffrey E. Hinton, Oriol Vinyals, and Jeffrey Dean. Distilling the knowledge in a neural network. *CoRR*, abs/1503.02531, 2015. 8, 13, 14, 29, 36, 39, 43, 51, 147
- [62] Gao Huang, Zhuang Liu, and Kilian Q. Weinberger. Densely connected convolutional networks. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2017. 1, 4, 99, 148
- [63] Frank Hutter, Holger H Hoos, and Kevin Leyton-Brown. Sequential model-based optimization for general algorithm configuration. In *International conference on learning and intelligent optimization*, pages 507–523. Springer, 2011. 4
- [64] Frank Hutter, Lars Kotthoff, and Joaquin Vanschoren, editors. *Automatic Machine Learning: Methods, Systems, Challenges*. Springer, 2018. In press, available at <http://automl.org/book>. 2, 3
- [65] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, pages 448–456, 2015. 12, 137, 138
- [66] Ozan Irsoy and Ethem Alpaydin. Autoencoder trees. In *Proceedings of The 7th Asian Conference on Machine Learning, ACML 2015, Hong Kong, November 20-22, 2015.*, pages 378–390, 2015. 16
- [67] M. Jaderberg, K. Simonyan, A. Vedaldi, and A. Zisserman. Reading text in the wild with convolutional neural networks. *arXiv preprint arXiv:1412.1842*, 2014. 72, 81, 82

- [68] M. Jaderberg, K. Simonyan, A. Vedaldi, and A. Zisserman. Synthetic data and artificial neural networks for natural scene text recognition. *arXiv preprint arXiv:1406.2227*, 2014. 63, 72, 82
- [69] Eric Jang, Shixiang Gu, and Ben Poole. Categorical Reparameterization with Gumbel-Softmax. In *ICLR*, 2017. 20, 128, 134, 139
- [70] Yangqing Jia, Joshua T. Abbott, Joseph L. Austerweil, Thomas L. Griffiths, and Trevor Darrell. Visual concept learning: Combining machine vision and bayesian generalization on concept hierarchies. In *Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems 2013. Proceedings of a meeting held December 5-8, 2013, Lake Tahoe, Nevada, United States.*, pages 1842–1850, 2013. 13
- [71] Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross Girshick, Sergio Guadarrama, and Trevor Darrell. Caffe: Convolutional architecture for fast feature embedding. *arXiv preprint arXiv:1408.5093*, 2014. 40, 43, 54
- [72] Lukasz Kaiser, Aidan N Gomez, and Francois Chollet. Depthwise separable convolutions for neural machine translation. *arXiv preprint arXiv:1706.03059*, 2017. 19, 137
- [73] Kirthevasan Kandasamy, Willie Neiswanger, Jeff Schneider, Barnabas Poczos, and Eric P Xing. Neural architecture search with bayesian optimisation and optimal transport. In *Advances in Neural Information Processing Systems*, pages 2016–2025, 2018. 4
- [74] Dimosthenis Karatzas, Faisal Shafait, Seiichi Uchida, Masakazu Iwamura, Lluís Gomez i Bigorda, Sergi Robles Mestre, Joan Mas, David Fernández Mota, Jon Almazán, and Lluís-Pere de las Heras. ICDAR 2013 robust reading competition. In *2013 12th International Conference on Document Analysis and Recognition, Washington, DC, USA, August 25-28, 2013*, pages 1484–1493, 2013. 82

- [75] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *ICLR*, 2015. 140
- [76] Aaron Klein, Stefan Falkner, Simon Bartels, Philipp Hennig, and Frank Hutter. Fast bayesian optimization of machine learning hyperparameters on large datasets. *arXiv preprint arXiv:1605.07079*, 2016. 5
- [77] Peter Kotschieder, Madalina Fiterau, Antonio Criminisi, and Samuel Rota Bulò. Deep neural decision forests. In *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence, IJCAI 2016, New York, NY, USA, 9-15 July 2016*, pages 4190–4194, 2016. 16
- [78] Alex Krizhevsky. Learning multiple layers of features from tiny images, 2009. Technical Report <https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf>. 18, 140, 144
- [79] Alex Krizhevsky. Learning multiple layers of features from tiny images, 2009. Technical Report <https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf>. 32, 63, 72, 109
- [80] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems 25, Lake Tahoe, Nevada, United States.*, pages 1106–1114, 2012. 1, 2, 12, 26, 33, 43, 63, 73, 99, 121, 123, 127
- [81] Yann LeCun, Corinna Cortes, and Christopher JC Burges. The mnist database of handwritten digits. 1998. 140, 141
- [82] Yann LeCun, John S. Denker, and Sara A. Solla. Optimal brain damage. In *Advances in Neural Information Processing Systems 2, [NIPS Conference, Denver, Colorado, USA, November 27-30, 1989]*, pages 598–605, 1989. 17

- [83] Yann LeCun, Fu Jie Huang, Leon Bottou, et al. Learning methods for generic object recognition with invariance to pose and lighting. In *CVPR (2)*, pages 97–104. Citeseer, 2004. 140, 142
- [84] Chen-Yu Lee, Patrick W. Gallagher, and Zhuowen Tu. Generalizing pooling functions in convolutional neural networks: Mixed, gated, and tree. In *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics, AISTATS 2016, Cadiz, Spain, May 9-11, 2016*, pages 464–472, 2016. 16
- [85] Chen-Yu Lee, Saining Xie, Patrick W. Gallagher, Zhengyou Zhang, and Zhuowen Tu. Deeply-supervised nets. In *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics, AISTATS 2015, San Diego, California, USA, May 9-12, 2015*, 2015. 12
- [86] Honglak Lee, Roger B. Grosse, Rajesh Ranganath, and Andrew Y. Ng. Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations. In *Proceedings of the 26th Annual International Conference on Machine Learning, ICML 2009, Montreal, Quebec, Canada, June 14-18, 2009*, pages 609–616, 2009. 12
- [87] Hongyang Li, Xiaoyang Guo, Bo Dai, Wanli Ouyang, and Xiaogang Wang. Neural network encapsulation. In *ECCV*, 2018. 144
- [88] Lin, Min, Qiang Chen, and Shuicheng Yan. Network in network. In *International Conference on Learning Representations, 2014 (arXiv:1409.1556)*., 2014. 41, 74, 75, 86, 137, 138
- [89] Baoyuan Liu, Fereshteh Sadeghi, Marshall F. Tappen, Ohad Shamir, and Ce Liu. Probabilistic label trees for efficient large scale image classification. In *2013 IEEE Conference on Computer Vision and Pattern Recognition, Portland, OR, USA, June 23-28, 2013*, pages 843–850, 2013. 13

- [90] Hanxiao Liu, Karen Simonyan, Oriol Vinyals, Chrisantha Fernando, and Koray Kavukcuoglu. Hierarchical representations for efficient architecture search. *arXiv preprint arXiv:1711.00436*, 2017. 16
- [91] Hanxiao Liu, Karen Simonyan, and Yiming Yang. Darts: Differentiable architecture search. *arXiv preprint arXiv:1806.09055*, 2018. 4, 5
- [92] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*, pages 3431–3440, 2015. 15
- [93] David G Lowe. Distinctive image features from scale-invariant keypoints. *International journal of computer vision*, 60(2):91–110, 2004. 1, 2
- [94] Simon M. Lucas, Alex Panaretos, Luis Sosa, Anthony Tang, Shirley Wong, and Robert Young. ICDAR 2003 robust reading competitions. In *7th International Conference on Document Analysis and Recognition (ICDAR 2003), 2-Volume Set, 3-6 August 2003, Edinburgh, Scotland, UK*, pages 682–687, 2003. 82
- [95] Andrew L Maas, Awni Y Hannun, and Andrew Y Ng. Rectifier nonlinearities improve neural network acoustic models. *Proc. ICML*, 30:1, 2013. 12, 99
- [96] Andrew L Maas, Awni Y Hannun, and Andrew Y Ng. Rectifier nonlinearities improve neural network acoustic models. In *Proc. icml*, volume 30, page 3, 2013. 137, 138
- [97] Chris J Maddison, Andriy Mnih, and Yee Whye Teh. The concrete distribution: A continuous relaxation of discrete random variables. *arXiv preprint arXiv:1611.00712*, 2016. 139
- [98] Marcin Marszalek and Cordelia Schmid. Constructing category hierarchies for visual recognition. In *Computer Vision - ECCV 2008, 10th European Conference on*

- Computer Vision, Marseille, France, October 12-18, 2008, Proceedings, Part IV*, pages 479–491, 2008. 13
- [99] Geoffrey F Miller, Peter M Todd, and Shailesh U Hegde. Designing neural networks using genetic algorithms. In *ICGA*, volume 89, pages 379–384, 1989. 4
- [100] Calvin Murdock, Zhen Li, Howard Zhou, and Tom Duerig. Blockout: Dynamic model selection for hierarchical deep networks. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 2583–2591, 2016. 17
- [101] Andrew Y. Ng, Michael I. Jordan, and Yair Weiss. On spectral clustering: Analysis and an algorithm. In *Advances in Neural Information Processing Systems 14 [Neural Information Processing Systems: Natural and Synthetic, NIPS 2001, December 3-8, 2001, Vancouver, British Columbia, Canada]*, pages 849–856, 2001. 36
- [102] Vicente Ordonez, Jia Deng, Yejin Choi, Alexander C. Berg, and Tamara L. Berg. From large scale image categorization to entry-level categories. In *IEEE International Conference on Computer Vision, ICCV 2013, Sydney, Australia, December 1-8, 2013*, pages 2768–2775, 2013. 13
- [103] Hieu Pham, Melody Y Guan, Barret Zoph, Quoc V Le, and Jeff Dean. Efficient neural architecture search via parameter sharing. *arXiv preprint arXiv:1802.03268*, 2018. 16
- [104] Sachin Ravi and Hugo Larochelle. Optimization as a model for few-shot learning. In *International Conference on Learning Representations (ICLR)*, 2017. 117
- [105] Esteban Real, Sherry Moore, Andrew Selle, Saurabh Saxena, Yutaka Leon Suematsu, Quoc V. Le, and Alex Kurakin. Large-scale evolution of image classifiers. *CoRR*, abs/1703.01041, 2017. 16
- [106] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 779–788, 2016. 1, 127

- [107] Sara Sabour, Nicholas Frosst, and Geoffrey E Hinton. Dynamic routing between capsules. In *Advances in neural information processing systems*, pages 3856–3866, 2017. 10, 18, 19, 127, 128, 129, 130, 131, 132, 133, 140, 149
- [108] Ruslan Salakhutdinov, Antonio Torralba, and Joshua B. Tenenbaum. Learning to share visual appearance for multiclass object detection. In *The 24th IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2011, Colorado Springs, CO, USA, 20-25 June 2011*, pages 1481–1488, 2011. 13
- [109] Tim Salimans, Jonathan Ho, Xi Chen, Szymon Sidor, and Ilya Sutskever. Evolution strategies as a scalable alternative to reinforcement learning. *arXiv preprint arXiv:1703.03864*, 2017. 16
- [110] Shreyas Saxena and Jakob Verbeek. Convolutional neural fabrics. In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pages 4053–4061, 2016. 18
- [111] Pierre Sermanet, David Eigen, Xiang Zhang, Michaël Mathieu, Rob Fergus, and Yann LeCun. Overfeat: Integrated recognition, localization and detection using convolutional networks. In *International Conference on Learning Representations (ICLR)*, 2013. 99
- [112] Bobak Shahriari, Kevin Swersky, Ziyu Wang, Ryan P Adams, and Nando De Freitas. Taking the human out of the loop: A review of bayesian optimization. *Proceedings of the IEEE*, 104(1):148–175, 2015. 4
- [113] Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. *arXiv preprint arXiv:1701.06538*, 2017. 18

- [114] Richard Shin, Charles Packer, and Dawn Song. Differentiable neural network architecture search. *International Conference on Learning Representations Workshop*, 2018. 4
- [115] L Sifre and S Mallat. Rigid-motion scattering for texture classification [ph. d. thesis]. *Ecole Polytechnique, CMAP*, 2014. 19, 137
- [116] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations (ICLR)*, 2015. 26, 35, 40, 63, 99, 127
- [117] Josef Sivic, Bryan C. Russell, Andrew Zisserman, William T. Freeman, and Alexei A. Efros. Unsupervised discovery of visual object class hierarchies. In *2008 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2008)*, 24-26 June 2008, Anchorage, Alaska, USA, 2008. 13
- [118] Jasper Snoek, Hugo Larochelle, and Ryan P. Adams. Practical bayesian optimization of machine learning algorithms. In *Advances in Neural Information Processing Systems 25, Lake Tahoe, Nevada, United States.*, pages 2960–2968, 2012. 16
- [119] Jasper Snoek, Oren Rippel, Kevin Swersky, Ryan Kiros, Nadathur Satish, Narayanan Sundaram, Md. Mostofa Ali Patwary, Prabhat, and Ryan P. Adams. Scalable bayesian optimization using deep neural networks. In *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, pages 2171–2180, 2015. 16, 38, 42
- [120] Nitish Srivastava, Geoffrey E. Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15(1):1929–1958, 2014. 15, 40
- [121] Nitish Srivastava and Ruslan Salakhutdinov. Discriminative transfer learning with tree-based priors. In *Advances in Neural Information Processing Systems 26: 27th Annual*

- Conference on Neural Information Processing Systems 2013. Proceedings of a meeting held December 5-8, 2013, Lake Tahoe, Nevada, United States.*, pages 2094–2102, 2013. 13
- [122] Kenneth O Stanley, David B D’Ambrosio, and Jason Gauci. A hypercube-based encoding for evolving large-scale neural networks. *Artificial life*, 15(2):185–212, 2009. 4
- [123] Kenneth O Stanley and Risto Miikkulainen. Evolving neural networks through augmenting topologies. *Evolutionary computation*, 10(2):99–127, 2002. 4
- [124] Felipe Petroski Such, Vashisht Madhavan, Edoardo Conti, Joel Lehman, Kenneth O Stanley, and Jeff Clune. Deep neuroevolution: genetic algorithms are a competitive alternative for training deep neural networks for reinforcement learning. *arXiv preprint arXiv:1712.06567*, 2017. 16
- [125] Kevin Swersky, David Duvenaud, Jasper Snoek, Frank Hutter, and Michael A Osborne. Raiders of the lost architecture: Kernels for bayesian optimization in conditional parameter spaces. *arXiv preprint arXiv:1409.4011*, 2014. 4
- [126] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott E. Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*, pages 1–9, 2015. 63, 99
- [127] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017. 19, 129, 132, 133
- [128] Andreas Veit and Serge Belongie. Convolutional networks with adaptive inference graphs. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 3–18, 2018. 20

- [129] Oriol Vinyals, Charles Blundell, Tim Lillicrap, Koray Kavukcuoglu, and Daan Wierstra. Matching networks for one shot learning. In *Advances in Neural Information Processing Systems 29, Barcelona, Spain*, pages 3630–3638, 2016. 109, 110, 117
- [130] Kai Wang, Boris Babenko, and Serge J. Belongie. End-to-end scene text recognition. In *IEEE International Conference on Computer Vision, ICCV 2011, Barcelona, Spain, November 6-13, 2011*, pages 1457–1464, 2011. 82
- [131] David Warde-Farley, Andrew Rabinovich, and Dragomir Anguelov. Self-informed neural network structure learning. *CoRR*, abs/1412.6563, 2014. 13, 14, 36, 39, 43, 51, 147
- [132] Daan Wierstra, Faustino J. Gomez, and Jürgen Schmidhuber. Modeling systems with internal state using evoluno. In *Genetic and Evolutionary Computation Conference, GECCO 2005, Proceedings, Washington DC, USA, June 25-29, 2005*, pages 1795–1802, 2005. 16
- [133] Felix Wu, Angela Fan, Alexei Baevski, Yann N Dauphin, and Michael Auli. Pay less attention with lightweight and dynamic convolutions. *ICLR*, 2019. 19, 137
- [134] Zuxuan Wu, Tushar Nagarajan, Abhishek Kumar, Steven Rennie, Larry S Davis, Kristen Grauman, and Rogerio Feris. Blockdrop: Dynamic inference paths in residual networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8817–8826, 2018. 18
- [135] Lingxi Xie and Alan L Yuille. Genetic cnn. In *ICCV*, pages 1388–1397, 2017. 16
- [136] Saining Xie, Ross B. Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2017. 4, 99, 102, 108, 109, 114, 116, 117, 118, 119, 123, 148
- [137] Zhicheng Yan, Hao Zhang, Robinson Piramuthu, Vignesh Jagadeesh, Dennis DeCoste, Wei Di, and Yizhou Yu. HD-CNN: hierarchical deep convolutional neural

- networks for large scale visual recognition. In *2015 IEEE International Conference on Computer Vision, ICCV 2015, Santiago, Chile, December 7-13, 2015*, pages 2740–2748, 2015. 13, 14, 29, 36, 37, 39, 41, 43, 51, 73, 74, 76, 147
- [138] Matthew D. Zeiler and Rob Fergus. Stochastic pooling for regularization of deep convolutional neural networks. *CoRR*, abs/1301.3557, 2013. 15
- [139] Matthew D. Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *Computer Vision - ECCV 2014 - 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part I*, pages 818–833, 2014. 99
- [140] Arber Zela, Aaron Klein, Stefan Falkner, and Frank Hutter. Towards automated deep learning: Efficient joint neural architecture and hyperparameter search. *arXiv preprint arXiv:1807.06906*, 2018. 5
- [141] Zhao Zhong, Junjie Yan, Wei Wu, Jing Shao, and Cheng-Lin Liu. Practical block-wise neural network architecture generation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2423–2432, 2018. 2, 4
- [142] Bolei Zhou, Àgata Lapedriza, Jianxiong Xiao, Antonio Torralba, and Aude Oliva. Learning deep features for scene recognition using places database. In *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*, pages 487–495, 2014. 63
- [143] Yang Zhou, Rong Jin, and Steven C. H. Hoi. Exclusive lasso for multi-task feature selection. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics, AISTATS 2010, Chia Laguna Resort, Sardinia, Italy, May 13-15, 2010*, pages 988–995, 2010. 30
- [144] Barret Zoph and Quoc V Le. Neural architecture search with reinforcement learning. *arXiv preprint arXiv:1611.01578*, 2016. 2, 4, 16

- [145] Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V Le. Learning transferable architectures for scalable image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8697–8710, 2018. 4, 5