

BIBLIOGRAPHY

- [1] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.
- [2] ———, “Identity mappings in deep residual networks,” in *European Conference on Computer Vision (ECCV)*. Springer, 2016, pp. 630–645.
- [3] A. Romero, N. Ballas, S. E. Kahou, A. Chassang, C. Gatta, and Y. Bengio, “Fitnets: Hints for thin deep nets,” in *Proceedings of International Conference on Learning Representations (ICLR)*, 2015.
- [4] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” *arXiv preprint arXiv:1409.1556*, 2014.
- [5] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, “Going deeper with convolutions,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 1–9.
- [6] S. Zagoruyko and N. Komodakis, “Wide residual networks,” in *Proceedings of the British Machine Vision Conference (BMVC)*, 2016, pp. 87.1–87.12.
- [7] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, “Aggregated residual transformations for deep neural networks,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 1492–1500.

- [8] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in Neural Information Processing Systems*, 2012, pp. 1097–1105.
- [9] D. Ciregan, U. Meier, and J. Schmidhuber, “Multi-column deep neural networks for image classification,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012, pp. 3642–3649.
- [10] M. Wang, “Multi-path convolutional neural networks for complex image classification,” *arXiv preprint arXiv:1506.04701*, 2015.
- [11] K. Kahatapitiya, D. Tissera, and R. Rodrigo, “Context-aware automatic occlusion removal,” in *2019 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2019, pp. 1895–1899.
- [12] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei, “ImageNet Large Scale Visual Recognition Challenge,” *International Journal of Computer Vision (IJCV)*, vol. 115, no. 3, pp. 211–252, 2015.
- [13] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.
- [14] D. Tissera, K. Kahatapitiya, R. Wijesinghe, S. Fernando, and R. Rodrigo, “Context-aware multipath networks,” *arXiv preprint arXiv:1907.11519*, 2019.
- [15] D. Tissera, K. Vithanage, R. Wijesinghe, K. Kahatapitiya, S. Fernando, and R. Rodrigo, “Feature-dependent cross-connections in multi-path neural networks,” in *2020 25th International Conference on Pattern Recognition (ICPR)*. IEEE, 2021, pp. 4032–4039.

- [16] D. Tissera, R. Wijesinghe, K. Vithanage, A. Xavier, S. Fernando, and R. Rodrigo, “End-to-end data-dependent routing in multi-path neural networks,” *Neural Computing and Applications*, pp. 1–20, 2023.
- [17] D. Tissera, K. Vithanage, R. Wijesinghe, A. Xavier, S. Jayasena, S. Fernando, and R. Rodrigo, “Neural mixture models with expectation-maximization for end-to-end deep clustering,” *Neurocomputing*, vol. 505, pp. 249–262, 2022.
- [18] K. Hornik, M. Stinchcombe, and H. White, “Multilayer feedforward networks are universal approximators,” *Neural networks*, vol. 2, no. 5, pp. 359–366, 1989.
- [19] P. Werbos, “Beyond regression: New tools for prediction and analysis in the behavioral sciences,” *PhD thesis, Committee on Applied Mathematics, Harvard University, Cambridge, MA*, 1974.
- [20] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, “Learning representations by back-propagating errors,” *Nature*, vol. 323, p. 533, 1986.
- [21] D. E. Rumelhart, G. E. Hinton, R. J. Williams *et al.*, “Learning internal representations by error propagation,” 1985.
- [22] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, “Gradient-based learning applied to document recognition,” *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [23] M. I. Jordan, “Serial order: A parallel distributed processing approach,” in *Advances in psychology*. Elsevier, 1997, vol. 121, pp. 471–495.
- [24] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [25] K. Cho, B. Van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, “Learning phrase representations using

- rnn encoder-decoder for statistical machine translation,” *arXiv preprint arXiv:1406.1078*, 2014.
- [26] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” *Advances in neural information processing systems*, vol. 30, 2017.
- [27] M. Zhou, Y. Bai, W. Zhang, T. Zhao, and T. Mei, “Look-into-object: Self-supervised structure modeling for object recognition,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 11 774–11 783.
- [28] X. Chen, C. Liang, D. Huang, E. Real, K. Wang, Y. Liu, H. Pham, X. Dong, T. Luong, C.-J. Hsieh *et al.*, “Symbolic discovery of optimization algorithms,” *arXiv preprint arXiv:2302.06675*, 2023.
- [29] P. Wang, S. Wang, J. Lin, S. Bai, X. Zhou, J. Zhou, X. Wang, and C. Zhou, “One-peace: Exploring one general representation model toward unlimited modalities,” *arXiv preprint arXiv:2305.11172*, 2023.
- [30] Y. Zhang, J. Qin, D. S. Park, W. Han, C.-C. Chiu, R. Pang, Q. V. Le, and Y. Wu, “Pushing the limits of semi-supervised learning for automatic speech recognition,” *arXiv preprint arXiv:2010.10504*, 2020.
- [31] P.-Y. Huang, V. Sharma, H. Xu, C. Ryali, H. Fan, Y. Li, S.-W. Li, G. Ghosh, J. Malik, and C. Feichtenhofer, “Mavil: Masked audio-video learners,” *arXiv preprint arXiv:2212.08071*, 2022.
- [32] T. Zhou, Z. Ma, Q. Wen, L. Sun, T. Yao, W. Yin, R. Jin *et al.*, “Film: Frequency improved legendre memory model for long-term time series forecasting,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 12 677–12 690, 2022.

- [33] A. Zeng, M. Chen, L. Zhang, and Q. Xu, “Are transformers effective for time series forecasting?” in *Proceedings of the AAAI conference on artificial intelligence*, vol. 37, no. 9, 2023, pp. 11 121–11 128.
- [34] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, “Language models are few-shot learners,” *Advances in neural information processing systems*, vol. 33, pp. 1877–1901, 2020.
- [35] OpenAI, “Gpt-4 technical report,” 2023.
- [36] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative adversarial nets,” in *Adv. in Neural Inf. Process. Syst.*, 2014, pp. 2672–2680.
- [37] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, “High-resolution image synthesis with latent diffusion models,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 10 684–10 695.
- [38] S. Gao, P. Zhou, M.-M. Cheng, and S. Yan, “Masked diffusion transformer is a strong image synthesizer,” *arXiv preprint arXiv:2303.14389*, 2023.
- [39] Y. LeCun, Y. Bengio, and G. Hinton, “Deep learning,” *nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [40] I. Goodfellow, Y. Bengio, and A. Courville, *Deep learning*. MIT press, 2016.
- [41] S. Hochreiter, “The vanishing gradient problem during learning recurrent neural nets and problem solutions,” *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, vol. 6, no. 02, pp. 107–116, 1998.
- [42] T. Dietterich, “Overfitting and undercomputing in machine learning,” *ACM computing surveys (CSUR)*, vol. 27, no. 3, pp. 326–327, 1995.

- [43] S. Ioffe and C. Szegedy, “Batch normalization: Accelerating deep network training by reducing internal covariate shift,” *arXiv preprint arXiv:1502.03167*, 2015.
- [44] X. Glorot and Y. Bengio, “Understanding the difficulty of training deep feedforward neural networks,” in *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings, 2010, pp. 249–256.
- [45] K. He, X. Zhang, S. Ren, and J. Sun, “Delving deep into rectifiers: Surpassing human-level performance on imagenet classification,” in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1026–1034.
- [46] A. Krogh and J. Hertz, “A simple weight decay can improve generalization,” *Advances in neural information processing systems*, vol. 4, 1991.
- [47] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, “Dropout: a simple way to prevent neural networks from overfitting,” *Journal of Machine Learning Research*, vol. 15, pp. 1929–1958, 2014.
- [48] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” in *Proceedings of International Conference on Learning Representations (ICLR)*, 2015.
- [49] S. Ruder, “An overview of gradient descent optimization algorithms,” *arXiv preprint arXiv:1609.04747*, 2016.
- [50] M. Tan and Q. Le, “Efficientnet: Rethinking model scaling for convolutional neural networks,” in *International conference on machine learning*. PMLR, 2019, pp. 6105–6114.
- [51] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. A. Alemi, “Inception-v4, inception-resnet and the impact of residual connections on learning,” in *AAAI Conference on Artificial Intelligence*, 2017.

- [52] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, “Rethinking the inception architecture for computer vision,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2818–2826.
- [53] Z. Chen, Y. Deng, Y. Wu, Q. Gu, and Y. Li, “Towards understanding mixture of experts in deep learning,” *arXiv preprint arXiv:2208.02813*, 2022.
- [54] R. Caruana, “Multitask learning,” *Machine learning*, vol. 28, no. 1, pp. 41–75, 1997.
- [55] K.-H. Thung and C.-Y. Wee, “A brief review on multi-task learning,” *Multimedia Tools and Applications*, vol. 77, no. 22, pp. 29 705–29 725, 2018.
- [56] M. Crawshaw, “Multi-task learning with deep neural networks: A survey,” *arXiv preprint arXiv:2009.09796*, 2020.
- [57] I. Misra, A. Shrivastava, A. Gupta, and M. Hebert, “Cross-stitch networks for multi-task learning,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016, pp. 3994–4003.
- [58] S. Ruder, J. Bingel, I. Augenstein, and A. Søgaard, “Latent multi-task architecture learning,” in *Proceedings of AAAI Conference of Artificial Intelligence*, February 2019, pp. 4822–4829.
- [59] Y. Gao, J. Ma, M. Zhao, W. Liu, and A. L. Yuille, “Nddr-cnn: Layerwise feature fusing in multi-task cnns by neural discriminative dimensionality reduction,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 3205–3214.
- [60] D. Ha, A. Dai, and Q. V. Le, “Hypernetworks,” in *Proceedings of International Conference on Learning Representations (ICLR)*, 2017.
- [61] G. E. Hinton, S. Sabour, and N. Frosst, “Matrix capsules with EM routing,” in *Proceedings of International Conference on Learning Representations (ICLR)*, 2018.

- [62] J. Hu, L. Shen, S. Albanie, G. Sun, and A. Vedaldi, “Gather-excite: Exploiting feature context in convolutional neural networks,” in *Advances in Neural Information Processing Systems*, 2018, pp. 9401–9411.
- [63] J. Hu, L. Shen, and G. Sun, “Squeeze-and-excitation networks,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018, pp. 7132–7141.
- [64] S. Sabour, N. Frosst, and G. E. Hinton, “Dynamic routing between capsules,” in *Advances in Neural Information Processing Systems*, 2017, pp. 3856–3866.
- [65] Q. Wang, B. Wu, P. Zhu, P. Li, W. Zuo, and Q. Hu, “Eca-net: Efficient channel attention for deep convolutional neural networks,” *arXiv preprint arXiv:1910.03151*, 2019.
- [66] A. Veit and S. Belongie, “Convolutional networks with adaptive inference graphs,” in *European Conference on Computer Vision*, 2018, pp. 3–18.
- [67] Z. Wu, T. Nagarajan, A. Kumar, S. Rennie, L. S. Davis, K. Grauman, and R. Feris, “Blockdrop: Dynamic inference paths in residual networks,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 8817–8826.
- [68] R. K. Srivastava, K. Greff, and J. Schmidhuber, “Highway networks,” *arXiv preprint arXiv:1505.00387*, 2015.
- [69] Y. Rao, J. Lu, J. Lin, and J. Zhou, “Runtime network routing for efficient image classification,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 41, no. 10, pp. 2291–2304, 2018.
- [70] X. Wang, F. Yu, Z.-Y. Dou, T. Darrell, and J. E. Gonzalez, “Skipnet: Learning dynamic routing in convolutional networks,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 409–424.

- [71] B. Chen, T. Zhao, J. Liu, and L. Lin, “Multipath feature recalibration densenet for image classification,” *International Journal of Machine Learning and Cybernetics*, vol. 12, no. 3, pp. 651–660, 2021.
- [72] H. Zhang, C. Wu, Z. Zhang, Y. Zhu, H. Lin, Z. Zhang, Y. Sun, T. He, J. Mueller, R. Manmatha *et al.*, “Resnest: Split-attention networks,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 2736–2746.
- [73] K. Yu, X. Wang, C. Dong, X. Tang, and C. C. Loy, “Path-restore: Learning network path selection for image restoration,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [74] R. K. Srivastava, K. Greff, and J. Schmidhuber, “Training very deep networks,” in *Advances in neural information processing systems*, 2015, pp. 2377–2385.
- [75] R. A. Jacobs, M. I. Jordan, S. J. Nowlan, and G. E. Hinton, “Adaptive mixtures of local experts,” *Neural computation*, vol. 3, no. 1, pp. 79–87, 1991.
- [76] M. I. Jordan and R. A. Jacobs, “Hierarchical mixtures of experts and the em algorithm,” *Neural computation*, vol. 6, no. 2, pp. 181–214, 1994.
- [77] D. Eigen, M. Ranzato, and I. Sutskever, “Learning factored representations in a deep mixture of experts,” *arXiv preprint arXiv:1312.4314*, 2013.
- [78] N. Shazeer, A. Mirhoseini, K. Maziarz, A. Davis, Q. Le, G. Hinton, and J. Dean, “Outrageously large neural networks: The sparsely-gated mixture-of-experts layer,” *arXiv preprint arXiv:1701.06538*, 2017.
- [79] W. Fedus, J. Dean, and B. Zoph, “A review of sparse expert models in deep learning,” *arXiv preprint arXiv:2209.01667*, 2022.
- [80] D. Lepikhin, H. Lee, Y. Xu, D. Chen, O. Firat, Y. Huang, M. Krikun, N. Shazeer, and Z. Chen, “Gshard: Scaling giant models with conditional

- computation and automatic sharding,” *arXiv preprint arXiv:2006.16668*, 2020.
- [81] W. Fedus, B. Zoph, and N. Shazeer, “Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity,” 2021.
- [82] C. Riquelme, J. Puigcerver, B. Mustafa, M. Neumann, R. Jenatton, A. Susano Pinto, D. Keysers, and N. Houlsby, “Scaling vision with sparse mixture of experts,” *Advances in Neural Information Processing Systems*, vol. 34, pp. 8583–8595, 2021.
- [83] L. Wu, M. Liu, Y. Chen, D. Chen, X. Dai, and L. Yuan, “Residual mixture of experts,” *arXiv preprint arXiv:2204.09636*, 2022.
- [84] A. Krizhevsky, G. Hinton *et al.*, “Learning multiple layers of features from tiny images,” Citeseer, Tech. Rep., 2009.
- [85] D. Ha, A. Dai, and Q. V. Le, “Hypernetworks,” *arXiv preprint arXiv:1609.09106*, 2016.
- [86] Facebook, “fb.resnet.torch.” [Online]. Available: <https://github.com/facebookarchive/fb.resnet.torch>
- [87] K. Simonyan, A. Vedaldi, and A. Zisserman, “Deep inside convolutional networks: Visualising image classification models and saliency maps,” *arXiv preprint arXiv:1312.6034*, 2013.
- [88] A. P. Dempster, N. M. Laird, and D. B. Rubin, “Maximum likelihood from incomplete data via the em algorithm,” *Journal of the Royal Statistical Society: Series B (Methodological)*, vol. 39, no. 1, pp. 1–22, 1977.
- [89] S. Laine and T. Aila, “Temporal ensembling for semi-supervised learning,” *arXiv preprint arXiv:1610.02242*, 2016.
- [90] X. Ji, J. F. Henriques, and A. Vedaldi, “Invariant information clustering for unsupervised image classification and segmentation,” in *Proceedings of*

- the International Conference on Computer Vision (ICCV)*, 2019, pp. 9865–9874.
- [91] B. Diallo, J. Hu, T. Li, G. A. Khan, X. Liang, and Y. Zhao, “Deep embedding clustering based on contractive autoencoder,” *Neurocomputing*, vol. 433, pp. 96–107, 2021.
- [92] S. Kullback and R. A. Leibler, “On information and sufficiency,” *The annals of mathematical statistics*, vol. 22, no. 1, pp. 79–86, 1951.
- [93] J. MacQueen *et al.*, “Some methods for classification and analysis of multivariate observations,” in *Proceedings of the fifth Berkeley Symposium on Mathematical Statistics and Probability*, vol. 1, no. 14. Oakland, CA, USA, 1967, pp. 281–297.
- [94] S. C. Johnson, “Hierarchical clustering schemes,” *Psychometrika*, vol. 32, no. 3, pp. 241–254, 1967.
- [95] G. J. McLachlan, S. X. Lee, and S. I. Rathnayake, “Finite mixture models,” *Annual review of statistics and its application*, vol. 6, pp. 355–378, 2019.
- [96] J. A. Hartigan, “Direct clustering of a data matrix,” *Journal of the American Statistical Association*, vol. 67, no. 337, pp. 123–129, 1972.
- [97] M. Ester, H.-P. Kriegel, J. Sander, X. Xu *et al.*, “A density-based algorithm for discovering clusters in large spatial databases with noise.” in *Kdd*, vol. 96, no. 34, 1996, pp. 226–231.
- [98] C. M. Bishop, “Pattern recognition and machine learning: springer new york,” 2006.
- [99] P. Haeusser, J. Plapp, V. Golkov, E. Aljalbout, and D. Cremers, “Associative deep clustering: Training a classification network with no labels,” in *German Conference on Pattern Recognition*. Springer, 2018, pp. 18–32.

- [100] M. Caron, P. Bojanowski, A. Joulin, and M. Douze, “Deep clustering for unsupervised learning of visual features,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 132–149.
- [101] J. Chang, L. Wang, G. Meng, S. Xiang, and C. Pan, “Deep adaptive image clustering,” in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 5879–5887.
- [102] J. Wu, K. Long, F. Wang, C. Qian, C. Li, Z. Lin, and H. Zha, “Deep comprehensive correlation mining for image clustering,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 8150–8159.
- [103] S. Han, S. Park, S. Park, S. Kim, and M. Cha, “Mitigating embedding and class assignment mismatch in unsupervised image classification,” in *16th European Conference on Computer Vision, ECCV 2020*. Springer, 2020.
- [104] J. Xie, R. Girshick, and A. Farhadi, “Unsupervised deep embedding for clustering analysis,” in *International Conference on Machine Learning (ICML)*, 2016, pp. 478–487.
- [105] X. Peng, I. W. Tsang, J. T. Zhou, and H. Zhu, “k-meansnet: When k-means meets differentiable programming,” *arXiv preprint arXiv:1808.07292*, 2018.
- [106] O. Kilinc and I. Uysal, “Learning latent representations in neural networks for clustering through pseudo supervision and graph-based activity regularization,” in *International Conference on Learning Representations*, 2018.
- [107] Y. Tao, K. Takagi, and K. Nakata, “Clustering-friendly representation learning via instance discrimination and feature decorrelation,” in *International Conference on Learning Representations*, 2020.
- [108] K. Greff, S. Van Steenkiste, and J. Schmidhuber, “Neural expectation maximization,” in *Advances in Neural Information Processing Systems*, 2017, pp. 6691–6701.

- [109] W. Van Gansbeke, S. Vandenhende, S. Georgoulis, M. Proesmans, and L. Van Gool, “Scan: Learning to classify images without labels,” in *European Conference on Computer Vision*. Springer, 2020, pp. 268–285.
- [110] T. W. Tsai, C. Li, and J. Zhu, “Mice: Mixture of contrastive experts for unsupervised image clustering,” in *International Conference on Learning Representations*, 2020.
- [111] R. E. Shiffler, “Maximum z scores and outliers,” *The American Statistician*, vol. 42, no. 1, pp. 79–80, 1988.
- [112] N. Metropolis and S. Ulam, “The monte carlo method,” *Journal of the American Statistical Association*, vol. 44, no. 247, pp. 335–341, 1949.
- [113] Y. LeCun, L. Bottou, Y. Bengio, P. Haffner *et al.*, “Gradient-based learning applied to document recognition,” *Proc. of the IEEE*, vol. 86, pp. 2278–2324, 1998.
- [114] S. Wold, K. Esbensen, and P. Geladi, “Principal component analysis,” *Chemometrics and intelligent laboratory systems*, vol. 2, no. 1-3, pp. 37–52, 1987.
- [115] L. v. d. Maaten and G. Hinton, “Visualizing data using t-sne,” *Journal of Machine Learning Research (MLR)*, vol. 9, pp. 2579–2605, 2008.
- [116] M. A. Kramer, “Nonlinear principal component analysis using autoassociative neural networks,” *AIChE journal*, vol. 37, no. 2, pp. 233–243, 1991.
- [117] A. Coates, A. Ng, and H. Lee, “An analysis of single-layer networks in unsupervised feature learning,” in *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, 2011, pp. 215–223.
- [118] L. Zelnik-Manor and P. Perona, “Self-tuning spectral clustering,” in *Advances in Neural Information Processing Systems*, 2005, pp. 1601–1608.

- [119] J. Wang, J. Wang, J. Song, X.-S. Xu, H. T. Shen, and S. Li, “Optimized cartesian k-means,” *IEEE Transactions on Knowledge and Data Engineering*, vol. 27, no. 1, pp. 180–192, 2014.
- [120] J. Yang, D. Parikh, and D. Batra, “Joint unsupervised learning of deep representations and image clusters,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 5147–5156.
- [121] M. Schultz and T. Joachims, “Learning a distance metric from relative comparisons,” in *Advances in Neural Information Processing Systems*, 2004, pp. 41–48.
- [122] Y. Bengio, P. Lamblin, D. Popovici, and H. Larochelle, “Greedy layer-wise training of deep networks,” in *Advances in Neural Information Processing Systems*, 2007, pp. 153–160.
- [123] A. Ng *et al.*, “Sparse autoencoder,” *CS294A Lecture notes*, vol. 72, no. 2011, pp. 1–19, 2011.
- [124] P. Vincent, H. Larochelle, I. Lajoie, Y. Bengio, P.-A. Manzagol, and L. Bottou, “Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion.” *Journal of Machine Learning Research (MLR)*, vol. 11, no. 12, 2010.
- [125] D. P. Kingma and M. Welling, “Auto-encoding variational bayes,” *arXiv preprint arXiv:1312.6114*, 2013.
- [126] J. Zhao, M. Mathieu, R. Goroshin, and Y. Lecun, “Stacked what-where auto-encoders,” *arXiv preprint arXiv:1506.02351*, 2015.
- [127] A. Radford, L. Metz, and S. Chintala, “Unsupervised representation learning with deep convolutional generative adversarial networks,” *arXiv preprint arXiv:1511.06434*, 2015.

- [128] A. Kosiorek, S. Sabour, Y. W. Teh, and G. E. Hinton, “Stacked capsule autoencoders,” in *Advances in Neural Information Processing Systems*, 2019, pp. 15 512–15 522.
- [129] H. W. Kuhn, “The hungarian method for the assignment problem,” *Naval research logistics quarterly*, vol. 2, no. 1-2, pp. 83–97, 1955.
- [130] S. Park, S. Han, S. Kim, D. Kim, S. Park, S. Hong, and M. Cha, “Improving unsupervised image clustering with robust learning,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 12 278–12 287.