# Learning to solve symbolic math from visual inputs

Thao Phung
University of Wyoming
phungpthao@gmail.com

Anh Nguyen
University of Wyoming
anh.ng8@gmail.com

Jeff Clune
Uber AI Labs, University of Wyoming
jeffclune@uwyo.edu

## Abstract

*Deep neural networks (DNNs) have obtained the state-of-the-art performance on many challenging tasks by learning from only visual inputs such as recognizing and detecting objects from images [8, 7, 6], playing Atari games from the visual screen [10], predicting sounds from silent videos [12] and even detecting cancers from medical scans [2]. The main reason behind DNNs' impressive performance is their ability to automatically extract a hierarchy of abstract features from the inputs that are useful to solve a given problem [4]. For example, to identify a person in a photo, DNNs extract simple features such as colors at lowest layers and edges and textures at higher layers, then corners and textures, and then more complicated features like chin and eyes of the person [4, 11].*

*Given their impressive pattern recognition capability, can DNNs learn to extract the meanings behind visual symbols? In this research, we explore approaches to train DNNs to solve addition and subtraction problems from images of equations (Fig. 1, right column). This is a challenging task that requires DNNs to learn the underlying values represented by these digit and operator symbols in order to produce the correct answers to the held-out problems (e.g. finding the answer 5 to the equation 2 + 3 = ?).*

*We show that by training a state-of-the-art DNN [7] on addition images (e.g. 1 + 2) together with their correct labels (here, 3), the DNN predicts correctly 98% of the time on new images of the equations it had already seen (Fig. 1a & b). However, the DNN performed poorly at 15% accuracy (slightly above random chance) on the equations that it had never seen before (Fig. 1c). This suggests that the DNN only memorizes the visual patterns of the equation images but does not realize the actual values of the numbers to perform addition. In this paper, we pose this open, challenging "symbolic math" task to the community and explore two strategies for solving it. Specifically, we teach DNNs to perform addition and subtraction via: (1) curriculum learning i.e. an organized training strategy in which we present a series of gradually increasingly complex problems [1] and (2) indirect learning i.e. training a separate DNN to paint the weights for the main DNN that takes in an equation im-*
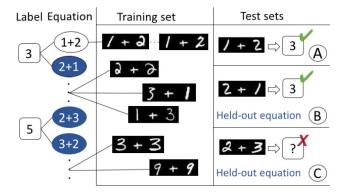
Figure 1: We train a DNN on images of addition problems (e.g. $1 + 2 =$?; note the training images only contain the left-hand side $1 + 2$ of a full equation) and ask the DNN to output the correct answer. Images are made up of two hand-written digits (taken from the MNIST dataset [9]) and a plus "+" sign. We produce 2000 images for each equation (top row, 3 example images of $1 + 2$), and include in the dataset all combinations of addition problems with two digits (e.g. $0 + 0, 0 + 1, ... 9 + 9$).

**(A)** As expected, the DNN performs well at recognizing new, unseen images of the $1 + 2$ problem that exists in the training set.

**(B)** Interestingly, the DNN can also correctly solve the $2 + 1 =$? problem when the training set only contains images of $1 + 2$ (not its commutative version $2 + 1$). This can be explained as the DNN learns to output "3" whenever it detects the visual patterns of "1", "+", and "2" in the image.

**(C)** However, when the commutative version is also held-out together with the original problem (e.g. both $2 + 3$ and $3 + 2$), the DNN is not able to correctly predict the answer. To solve this problem, the DNN would have to "understand" the underlying values of digits and symbols. We are exploring different strategies of harnessing curriculum learning [1] and indirect learning [5, 13, 3] to solve this challenging task.

*age as input and outputs the predicted answer [3, 13, 5]. We will report and discuss our results in the poster.*

# References

[1] Y. Bengio, J. Louradour, R. Collobert, and J. Weston. Curriculum learning. In *Proceedings of the 26th annual international conference on machine learning*, pages 41–48. ACM, 2009. 1

[2] D. C. Cireşan, A. Giusti, L. M. Gambardella, and J. Schmidhuber. Mitosis detection in breast cancer histology images with deep neural networks. In *International Conference on Medical Image Computing and Computer-assisted Intervention*, pages 411–418. Springer, 2013. 1

[3] C. Fernando, D. Banarse, M. Reynolds, F. Besse, D. Pfau, M. Jaderberg, M. Lanctot, and D. Wierstra. Convolution by evolution: Differentiable pattern producing networks. In *Proceedings of the 2016 on Genetic and Evolutionary Computation Conference*, pages 109–116. ACM, 2016. 1

[4] I. Goodfellow, Y. Bengio, and A. Courville. *Deep learning*. MIT Press, 2016. 1

[5] D. Ha, A. Dai, and Q. V. Le. Hypernetworks. arxiv preprint. *arXiv preprint arXiv:1609.09106*, 2016. 1

[6] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask r-cnn. *arXiv preprint arXiv:1703.06870*, 2017. 1

[7] 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*, pages 770–778, 2016. 1

[8] G. Huang, Z. Liu, K. Q. Weinberger, and L. van der Maaten. Densely connected convolutional networks. *arXiv preprint arXiv:1608.06993*, 2016. 1

[9] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. 1

[10] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015. 1

[11] A. Nguyen, J. Yosinski, and J. Clune. Multifaceted feature visualization: Uncovering the different types of features learned by each neuron in deep neural networks. *arXiv preprint arXiv:1602.03616*, 2016. 1

[12] A. Owens, P. Isola, J. McDermott, A. Torralba, E. H. Adelson, and W. T. Freeman. Visually indicated sounds. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2405–2413, 2016. 1

[13] K. O. Stanley. Compositional pattern producing networks: A novel abstraction of development. *Genetic programming and evolvable machines*, 8(2):131–162, 2007. 1