

The disadvantage of CNN versus DBN image classification under adversarial conditions

Tao (Andy) Yang†, Daniel L. Silver †,*
† Acadia University

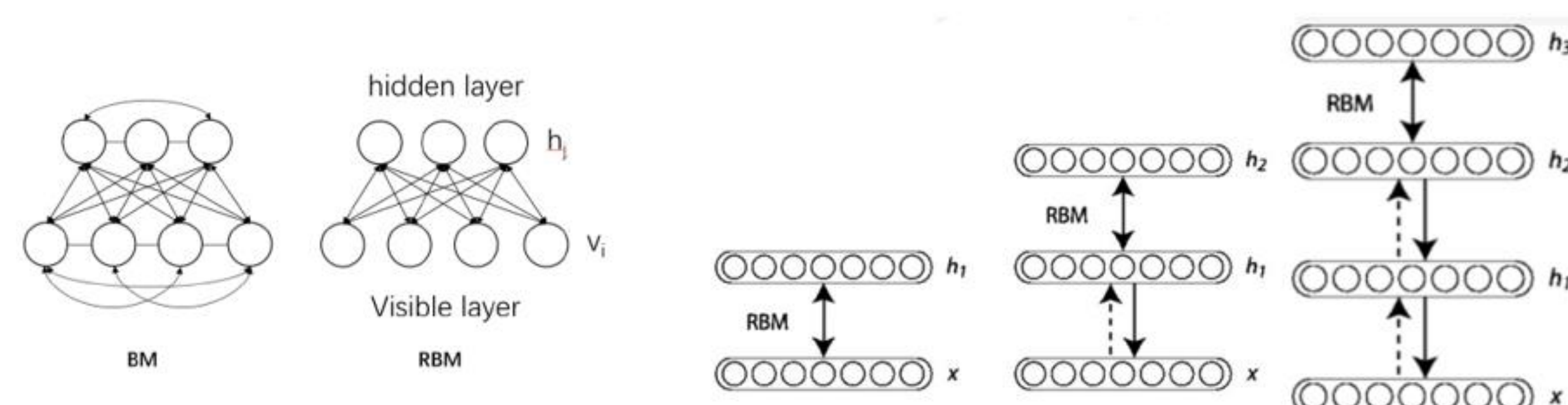
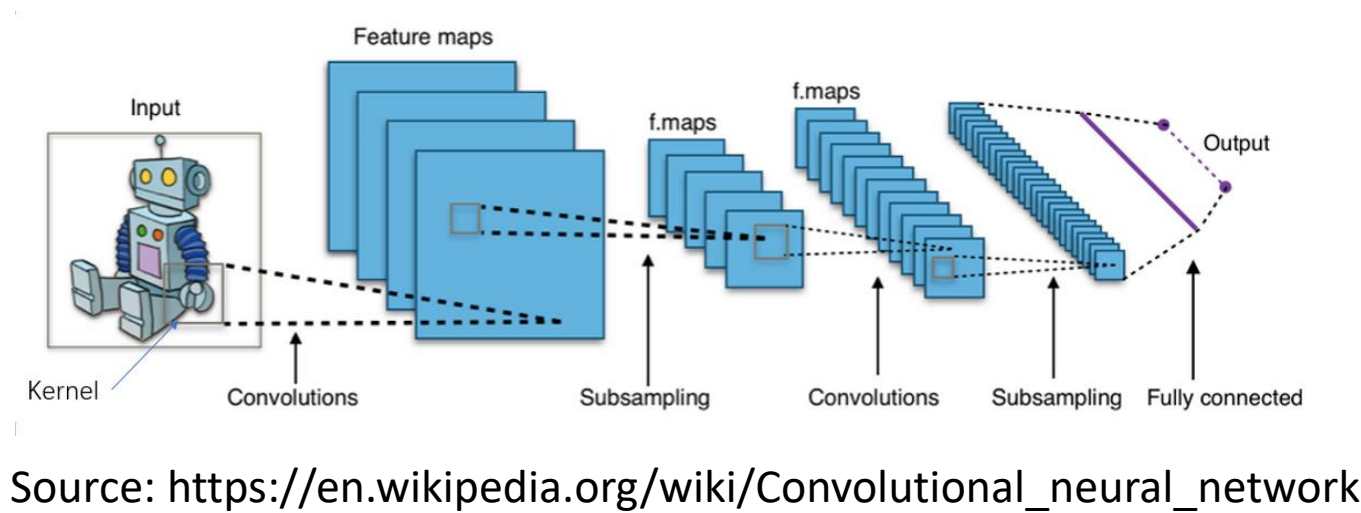
Summary

- We compare Convolutional Neural Networks (CNN) and Deep Belief Network's (DBN) ability to withstand common image classification attacks.
- The results show that the DBN models generally perform better under attack than the CNN models.

Introduction

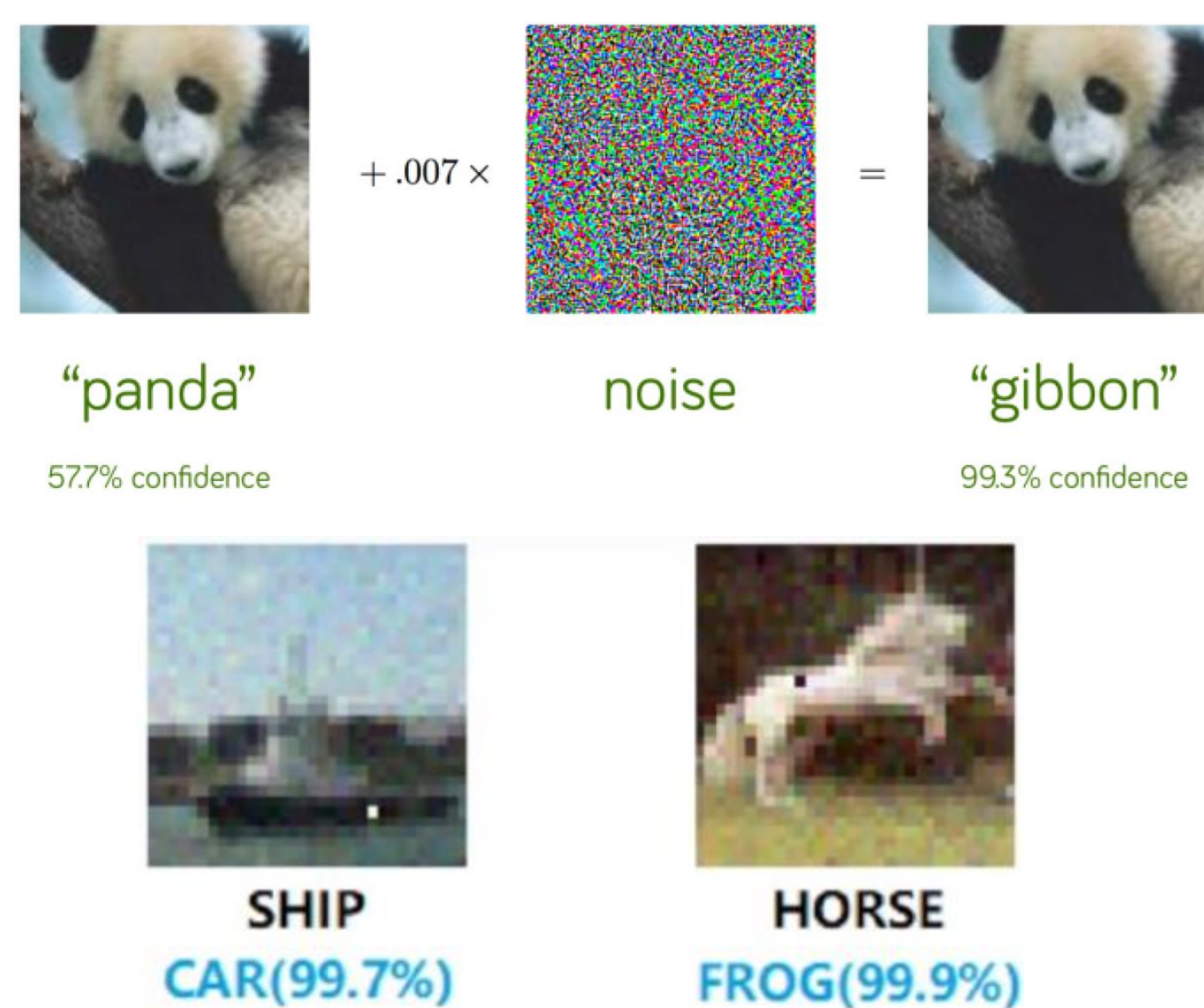
- Focus is on images created to attack the performance of CNN and DBN classification models. We show that the CNN inductive bias, which assumes that important features can be extracted from adjacent pixels, fails under certain adversarial attacks.

- In contrast, DBNs, which are pre-trained using unsupervised examples, learn features with no assumption about the proximity of pixels, and do better under the same attacks.



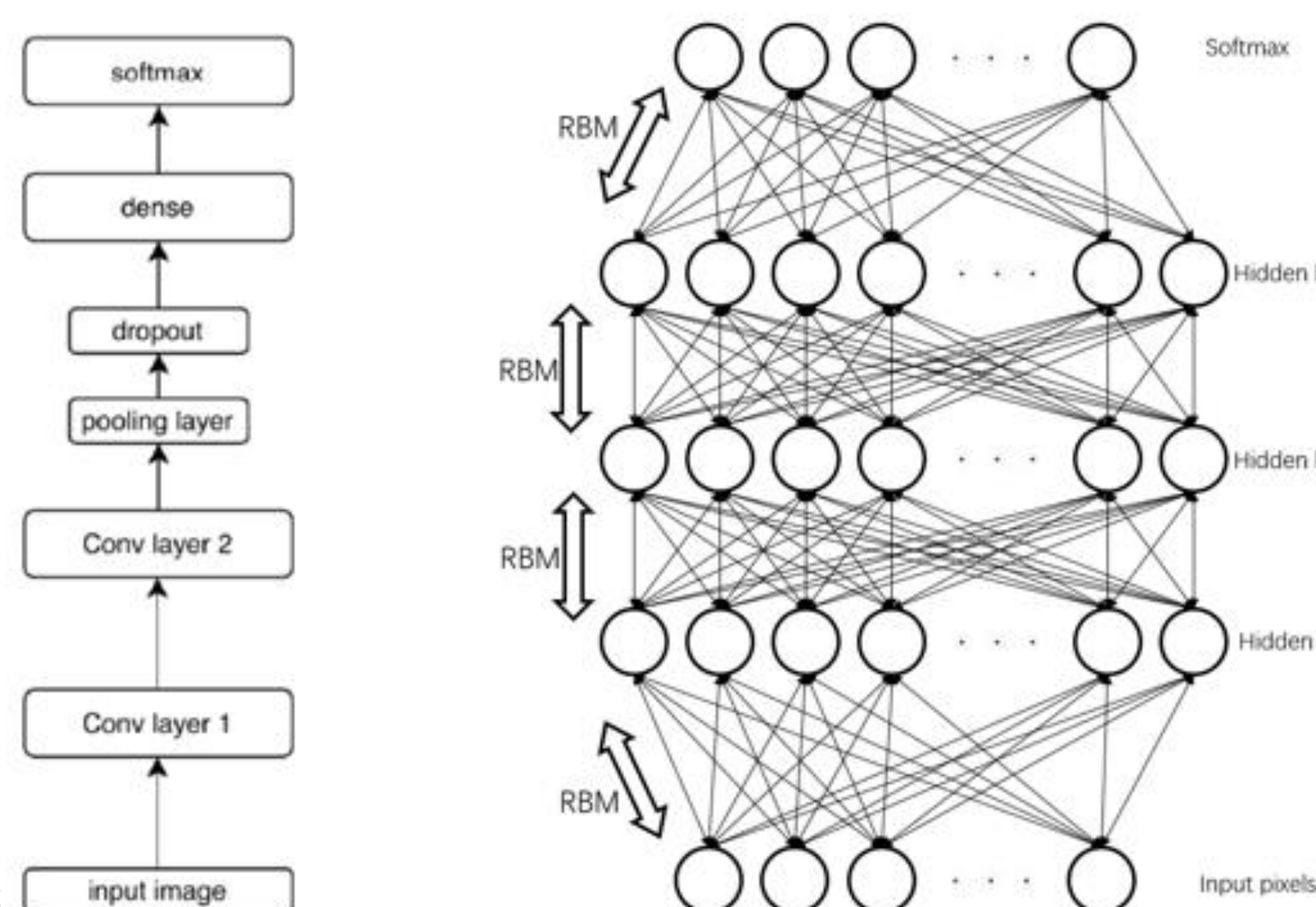
Background

- It is possible to trick trained neural networks by carefully modifying the pixels of the images.
- This is called an adversarial attack, and if done well, it will make the network think an image of a panda is an image of a gibbon.
- In 2018, it was shown that by only changing one pixel, a trained CNN network will incorrectly think a horse is a frog¹.
- Most images in the CIFAR-10 dataset can be modified similarly to fool CNN models



Theory and approach

- Our theory is that CNN makes a strong inductive bias assumption about the relationship between pixels that are proximal to each other, and this will be a disadvantage for CNN when people intentionally modify the pixels.
- We wish to show that the CNN inductive bias, according to proximity of pixels, which assumes pixels adjacent to each other will tend to have the same colour and brightness, fails under certain adversarial attacks.
- DBNs, which learn features with a less specific inductive bias with no assumptions about the proximity of pixels, will not fail in the same way².

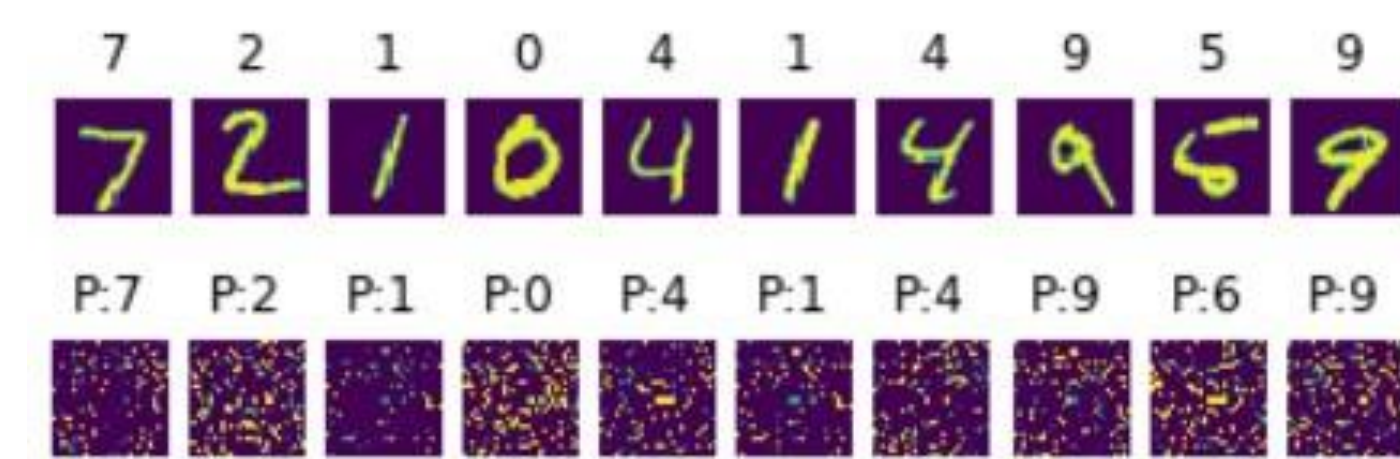


Acknowledgements

First and foremost, I would like to show my deepest gratitude to my supervisor, Dr. Daniel Silver. Without his help, the goal of this work would not have been realized. I shall extend my thanks to all professors in the computer science department at Acadia University for their help. I would also like to give a special thank to Dr. Jim Diamond, Dr. Greg Lee, Dr. Zhang Haiyi and Professor Duane Currie, as they gave me accurate and efficient help on many deep questions. They also provide extra convenience and help which is much more than I should get. Last but not least, I would like to thank all my friends and family for their financial and moral support.

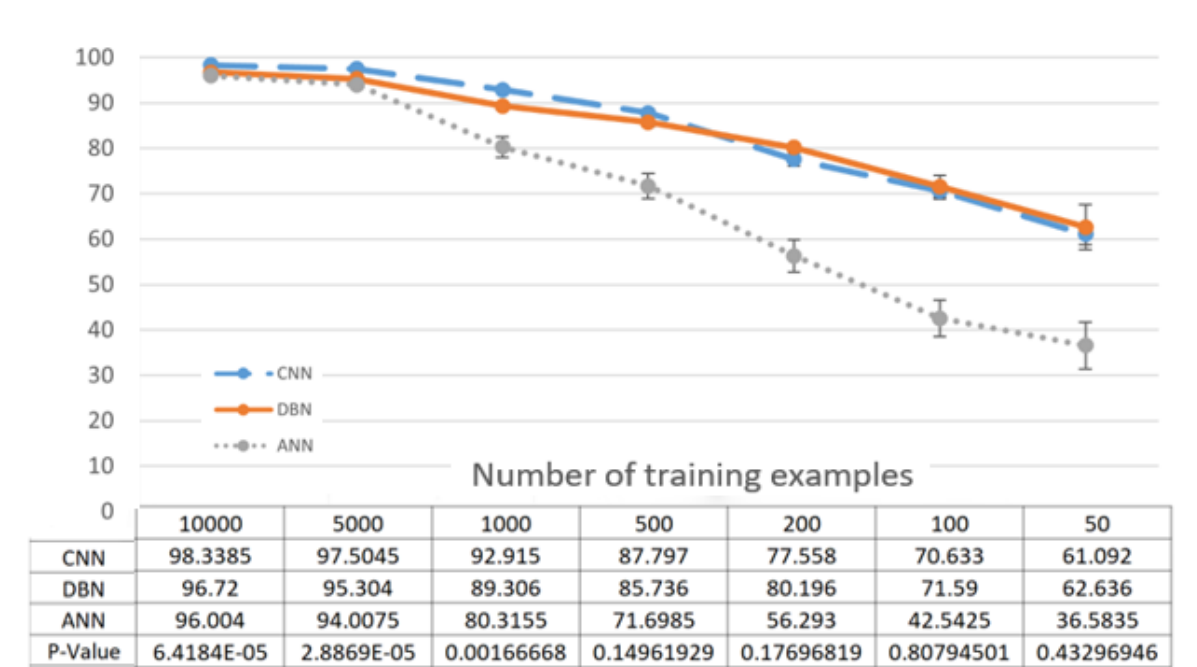
Empirical studies

- This section aims to demonstrate empirically that a semi-supervised approach using DBNs is superior to the strictly supervised CNN approach when under attack.
- The first experiment uses the MNIST data set to show the weakness of the CNN algorithm, as compared to the DBN algorithm, when the kernel assumption no longer holds for the input data.



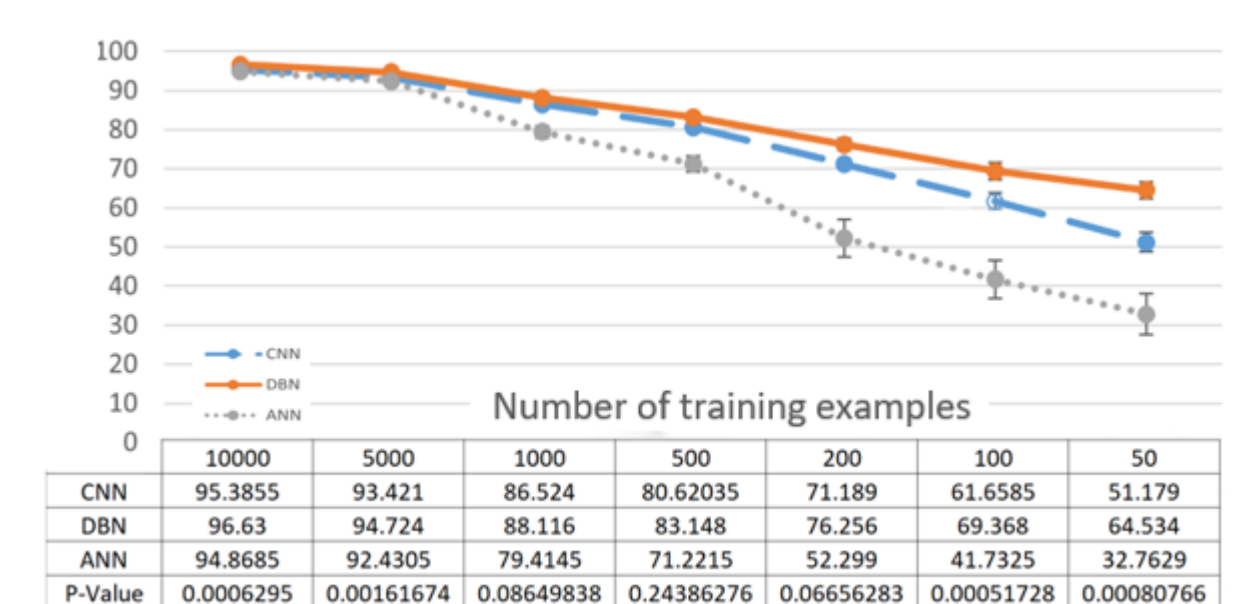
Examples of images used in the first experiment: top row shows original images, bottom row shows mixed-pixel images.

- CNN has an advantage over DBN when the training samples are high, but then loses the advantage as the numbers of labelled training examples decreases.
- This is expected since CNN relies on a large number of human-labeled training samples to learn an accurate model.



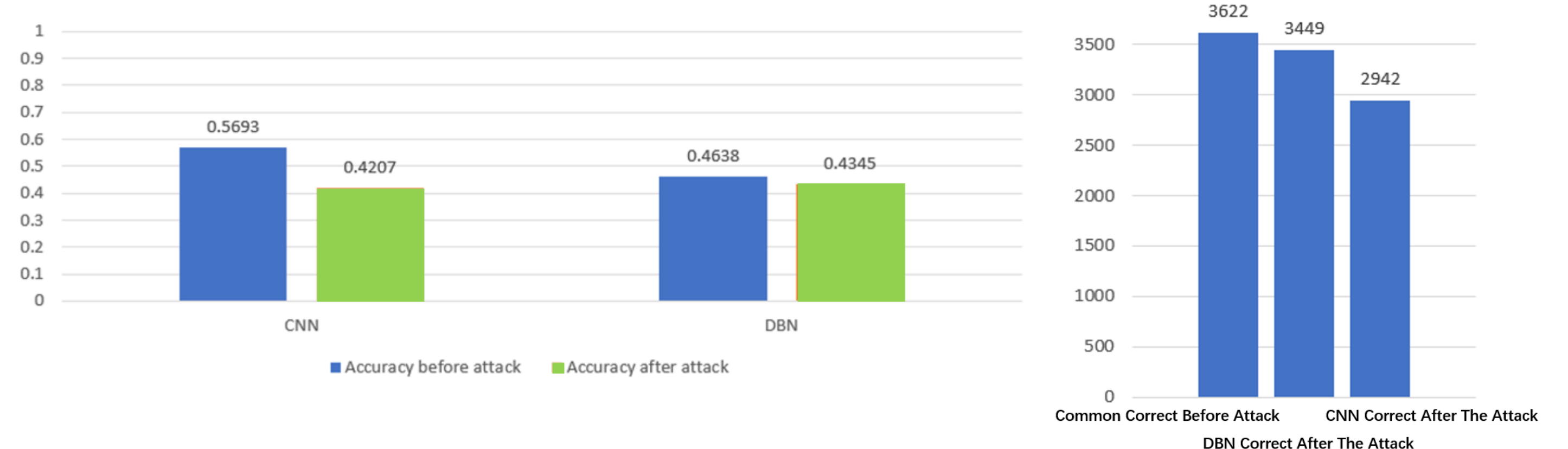
CNN and DBN accuracy on original test images by number of training examples.

- Mixing up the pixels of the images adversely affects the performance of the CNN models
- But DBN models are not affected – their level of accuracy is as before the attack



CNN and DBN accuracy on mixed-pixel test images by number of training examples.

- The second experiment tests the performance of CNN and DBN models against noisy one-pixel attack images and shows the significant loss in model accuracy that CNN suffers versus DBN.



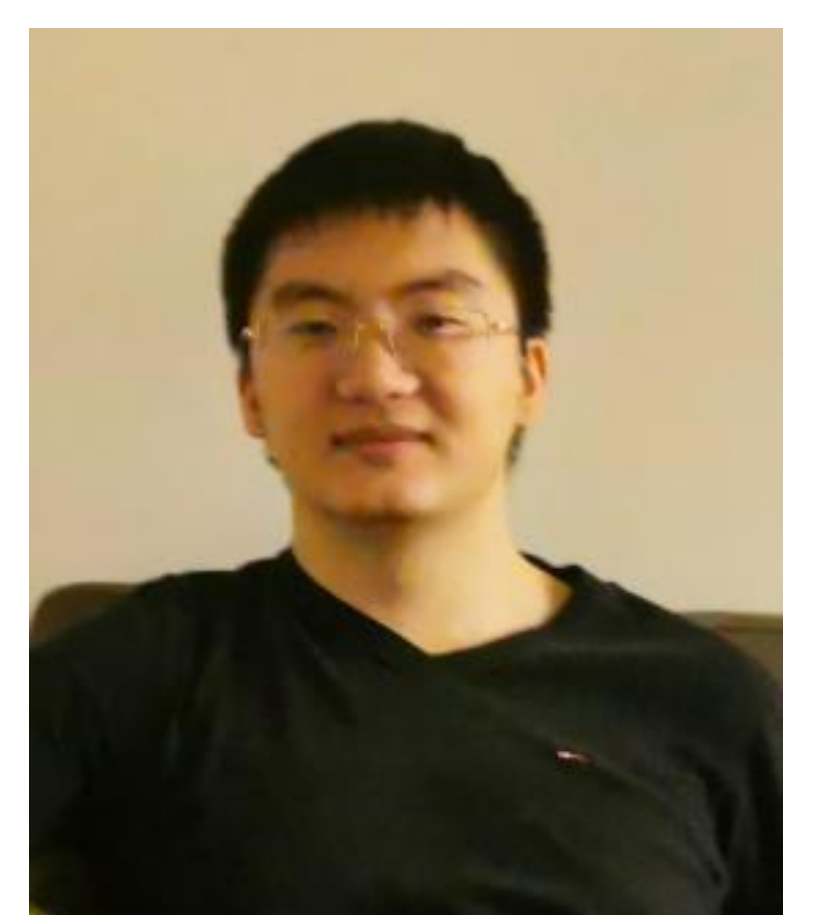
Left: Performance on test examples before and after the one-pixel attack. Right: Performance on test images before and after the one-pixel attack.

Conclusion

- CNNs inductive bias becomes a disadvantage when the assumption of the relationship between proximal pixels no longer holds.
- In comparison, the DBN algorithm assumes that the best internal representation can be developed by pre-training the network using large sets of unlabelled examples from the same input space without other a priori assumptions.
- This suggests that in the context of lifelong machine learning (or continual learning) systems that the DBN approach of learning the underlying internal representation is superior to the highly engineered CNN approach³.

References

- J. Su, D. V. Vargas, and K. Sakurai. "One Pixel Attack for Fooling Deep Neural Networks". In: IEEE Trans. on Evol. Comp. 23.5 (2019), pp. 828–841. doi: 10.1109/tevc.2019.2890858.
- G. E. Hinton. "A Practical Guide to Training Restricted Boltzmann Machines". In: Neural Networks: Tricks of the Trade, 2nd Edition, LNCS 7700 (2012), pp. 599–619.
- A. Anand. Contrastive Self-Supervised Learning. <https://ankeshanand.com/blog/2020/01/26/contrastive-self-supervised-learning.html>. 2020.



Tao Yang
137660y@acadiau.ca