

FUNDAMENTALS OF DEEP LEARNING FOR MULTI-GPUS

LAB 3, PART 1: SCALING THE BATCH SIZE



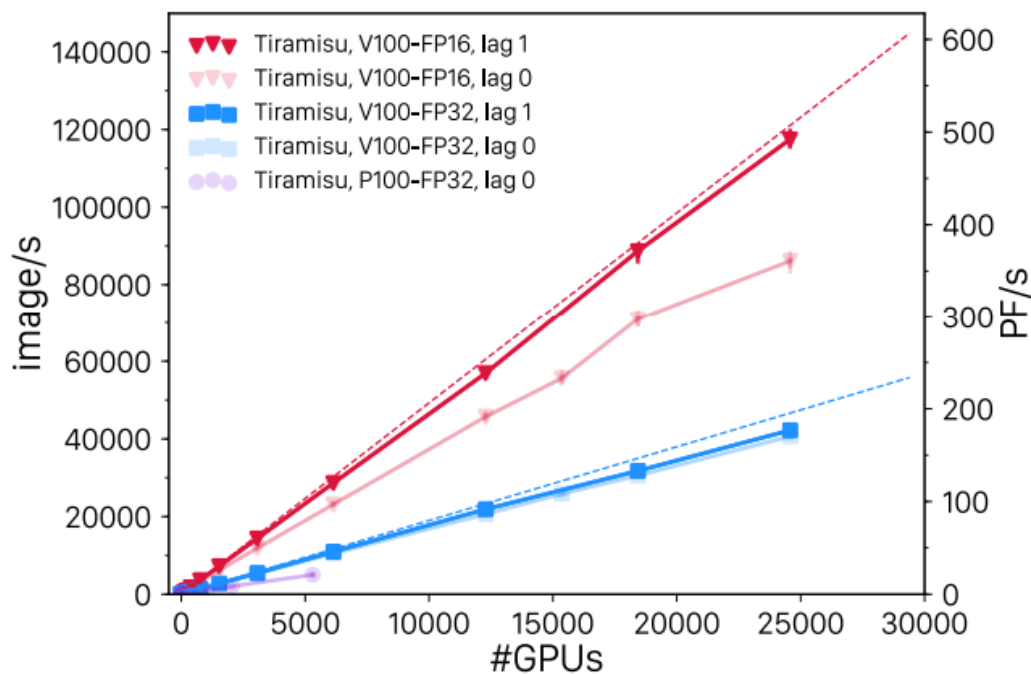
DEEP
LEARNING
INSTITUTE



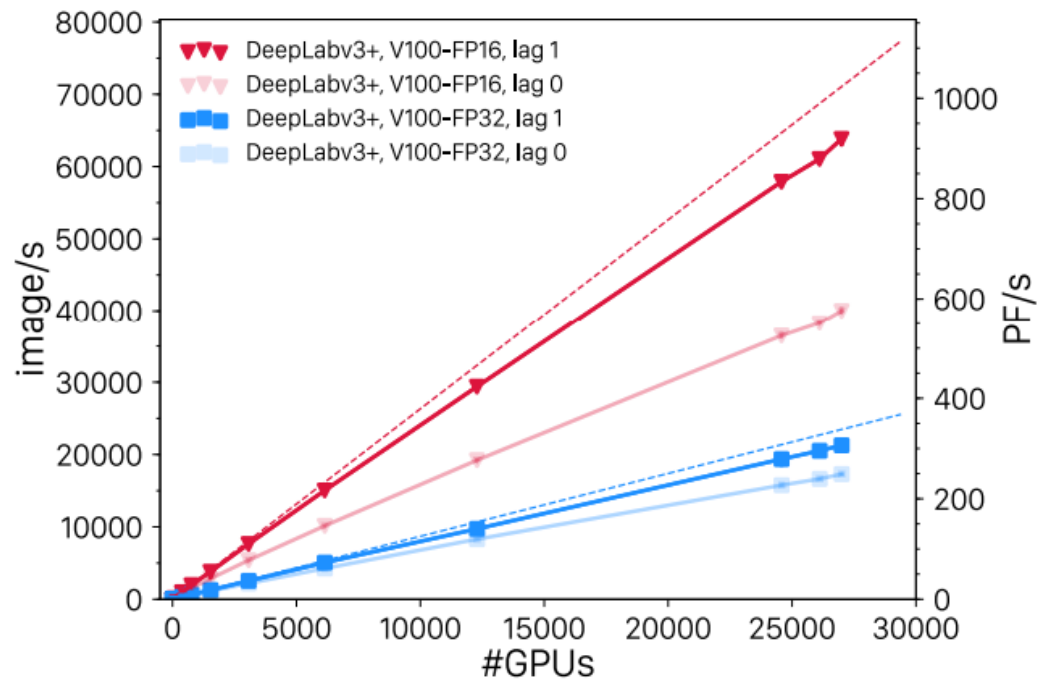
**CAN WE INCREASE THE BATCH SIZE
INDEFINITELY?**

IN TERMS OF IMAGES / SECOND?

Yes



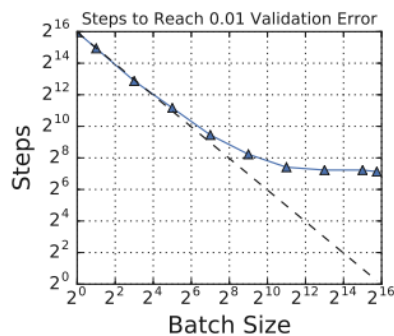
(a) Tiramisu



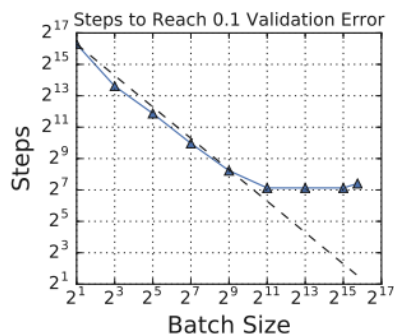
(b) DeepLabv3+

IN TERMS OF STEPS TO CONVERGENCE?

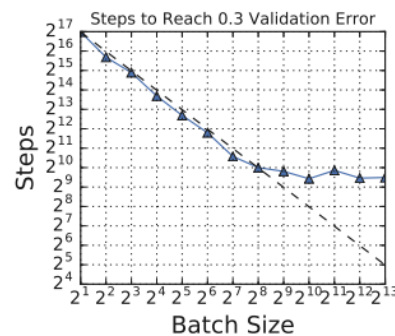
There are limits



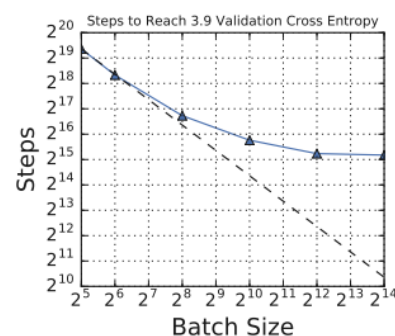
(a) Simple CNN on MNIST



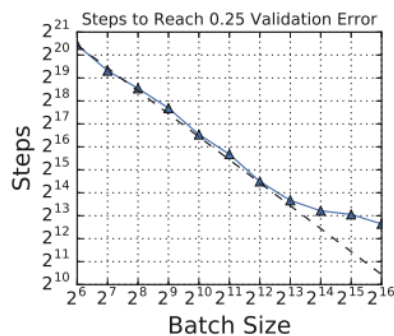
(b) Simple CNN on Fashion MNIST



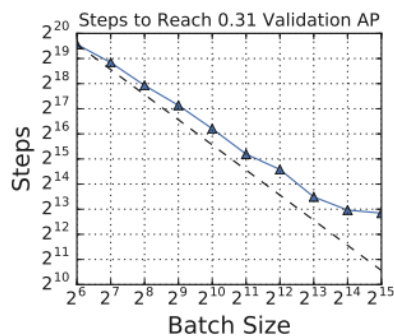
(c) ResNet-8 on CIFAR-10



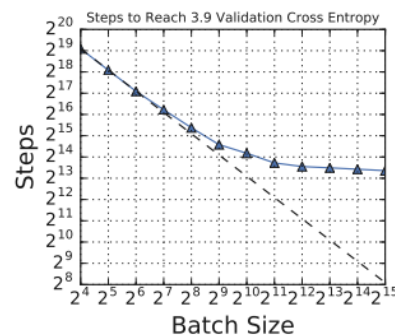
(g) Transformer on Common Crawl



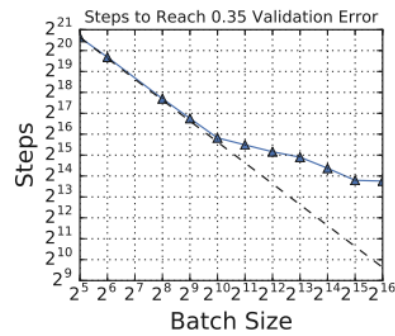
(d) ResNet-50 on ImageNet



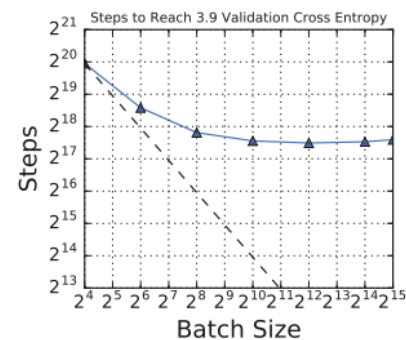
(e) ResNet-50 on Open Images



(f) Transformer on LM1B



(h) VGG-11 on ImageNet



(i) LSTM on LM1B

IN TERMS OF STEPS TO CONVERGENCE?

There are limits

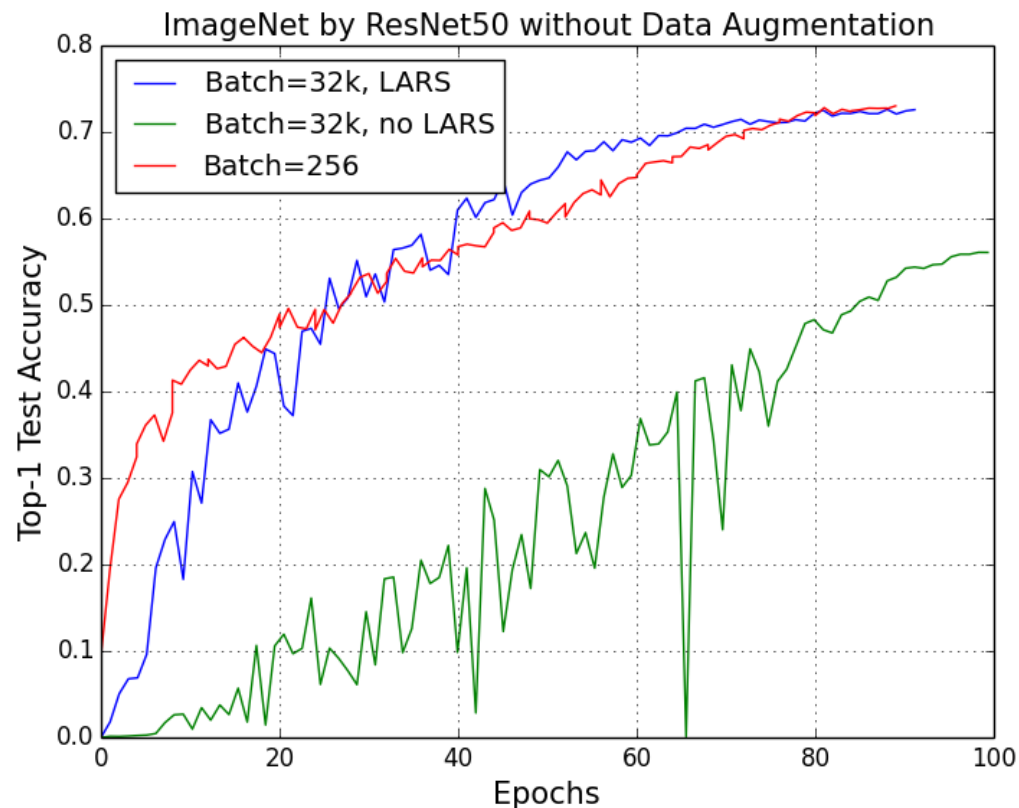
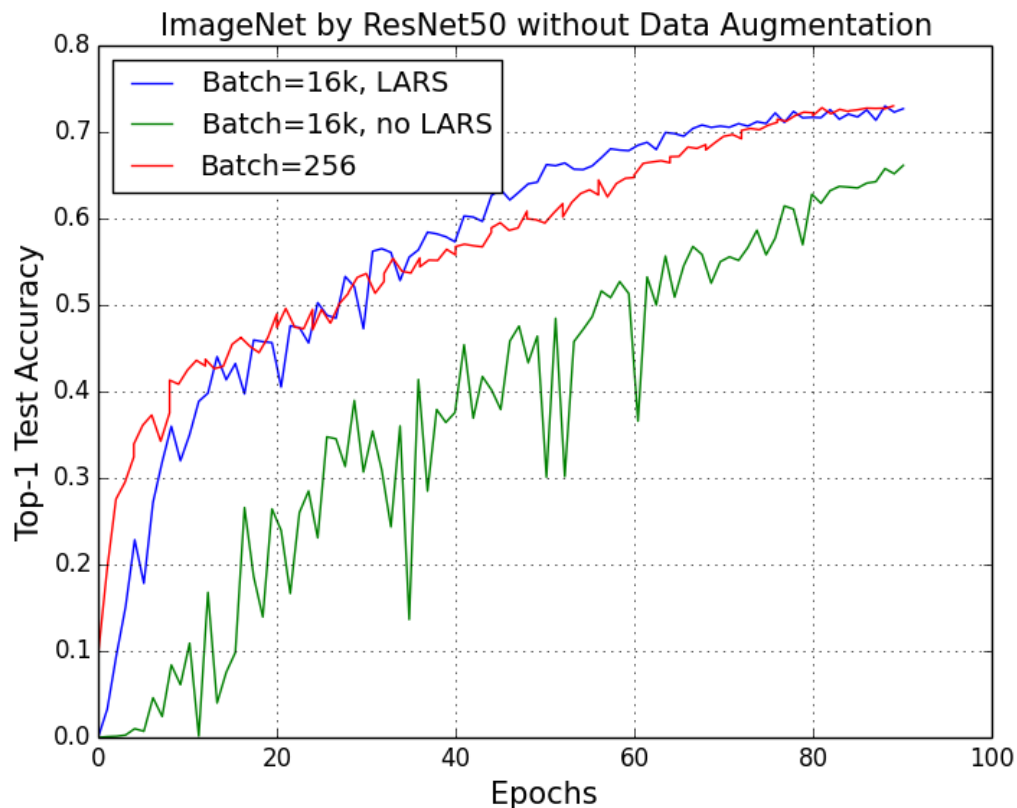




LARGE MINIBATCH AND ITS IMPACT ON ACCURACY

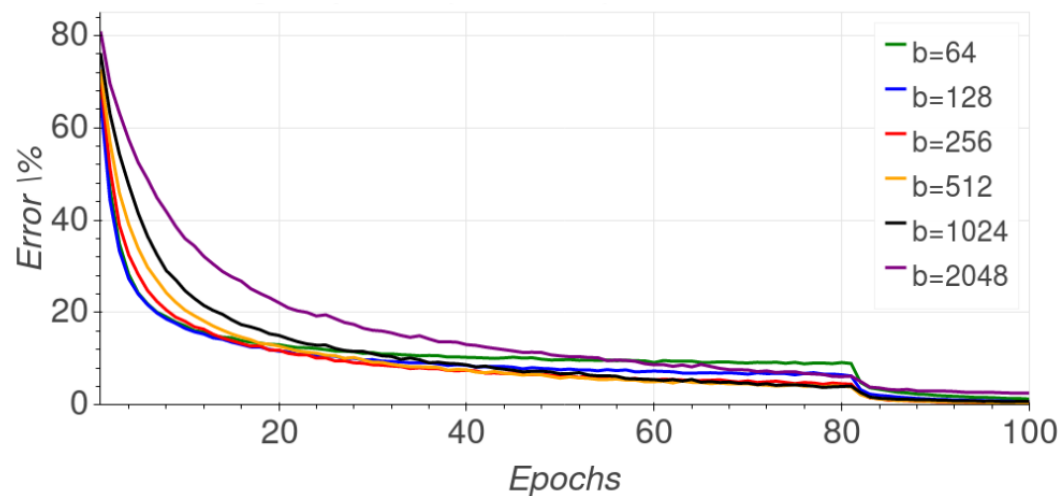
IMPACT ON ACCURACY

Naïve approaches lead to degraded accuracy

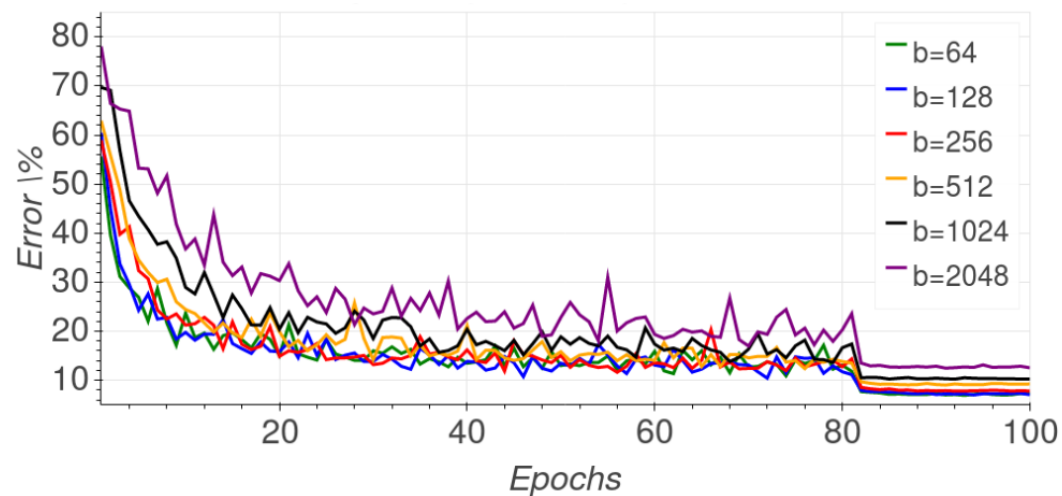


IMPACT ON ACCURACY

Naïve approaches lead to degraded accuracy



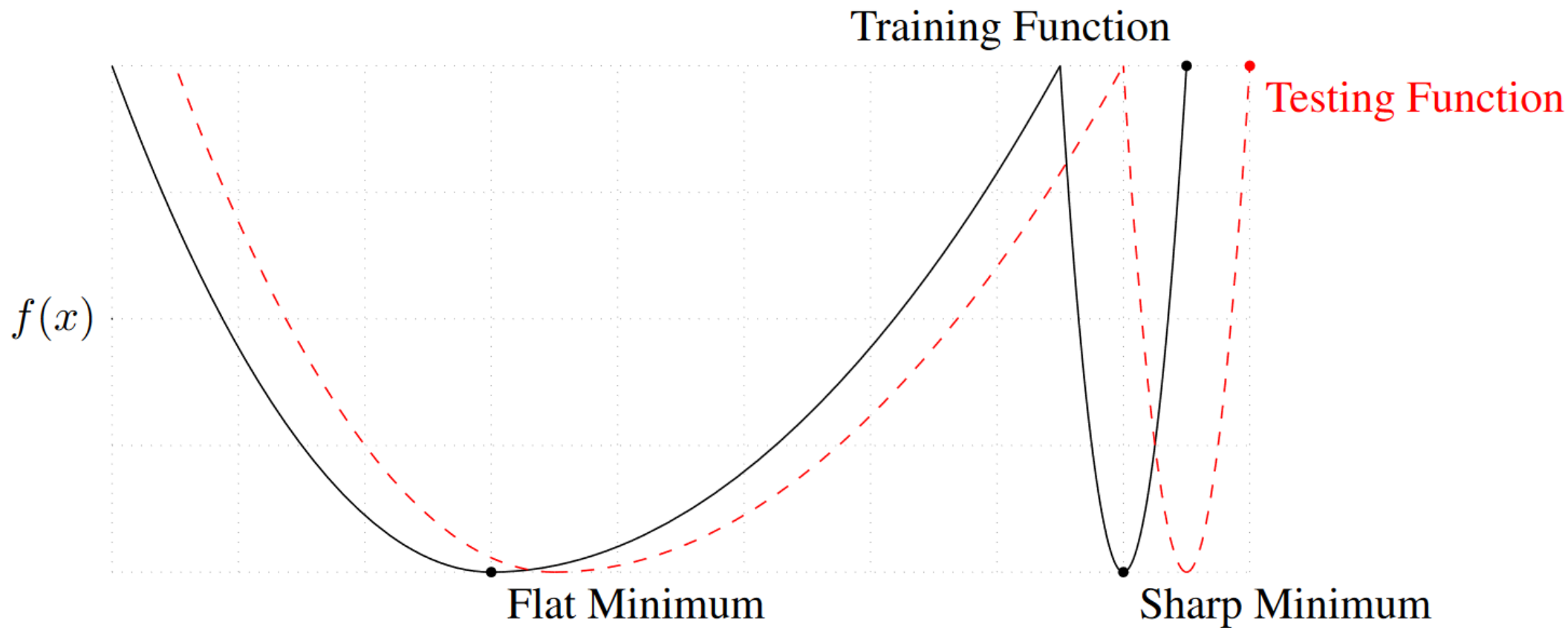
(a) Training error



(b) Validation error

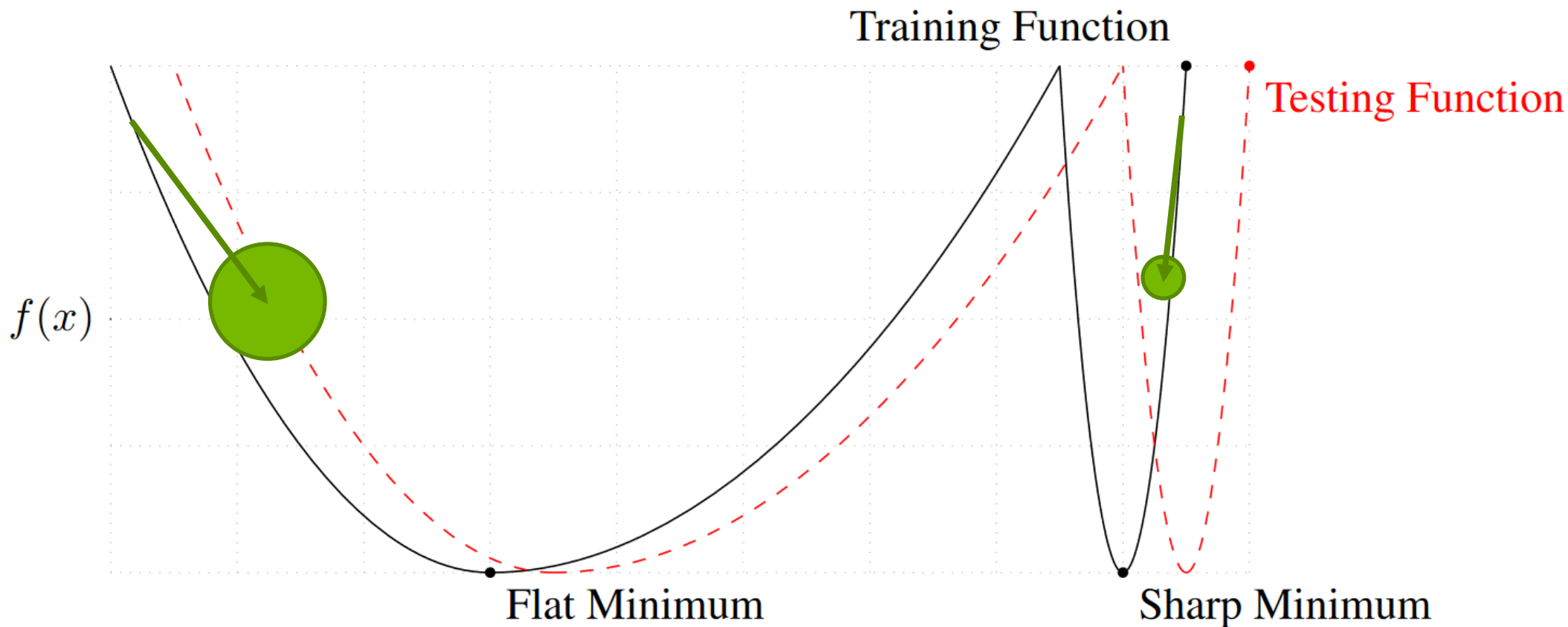
IMPACT ON ACCURACY

Why? Generalization and flatness of minima?



IMPACT ON ACCURACY

Why does it happen? Noise in the gradient update.



IMPACT ON ACCURACY

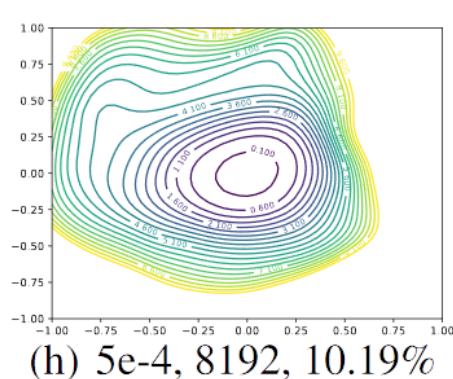
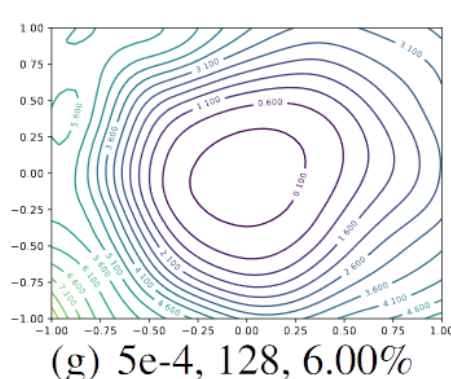
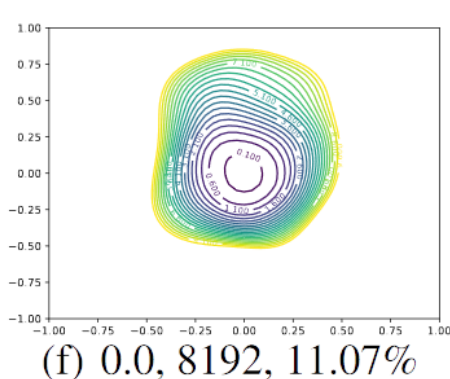
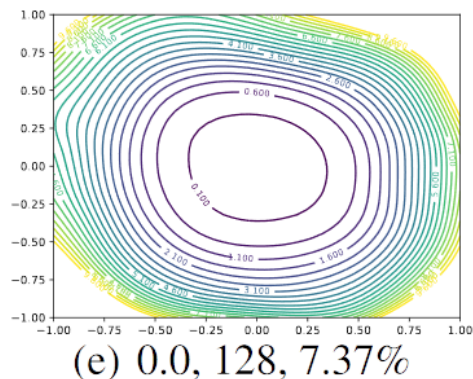
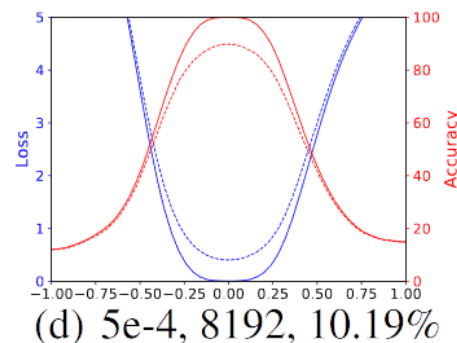
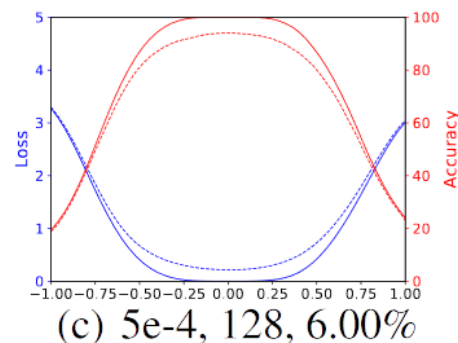
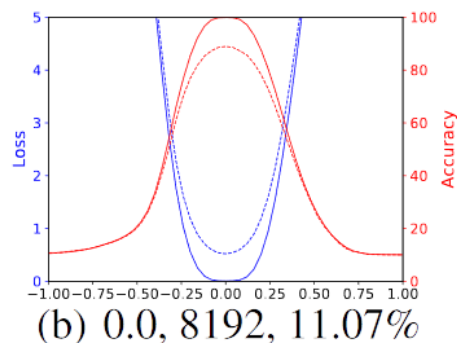
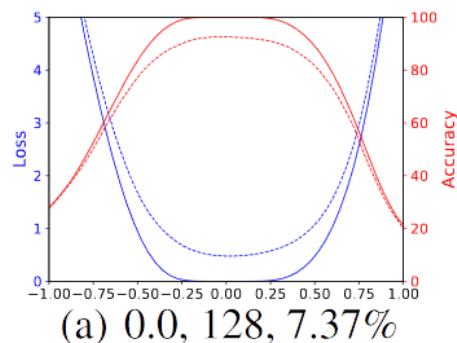
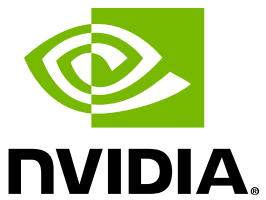


Figure 3: The 1D and 2D visualization of solutions obtained using SGD with different weight decay and batch size. The title of each subfigure contains the weight decay, batch size, and test error.



DEEP
LEARNING
INSTITUTE

www.nvidia.com/dli