**Applied Accelerated Artificial Intelligence**
**Prof. Bharathkumar**
**Department of Computer Science and Engineering**
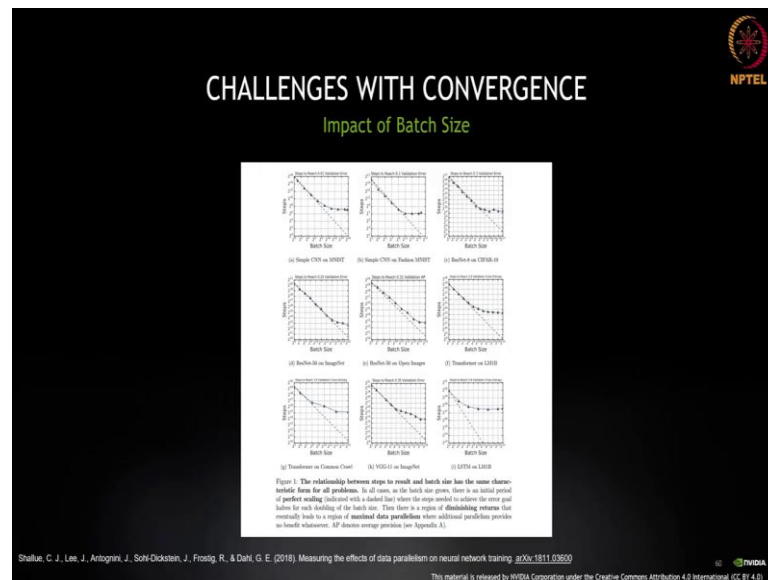**Indian Institute of Technology, Palakkad**

**Lecture - 39**
**Challenges with Distributed Deep Learning Training Convergence**
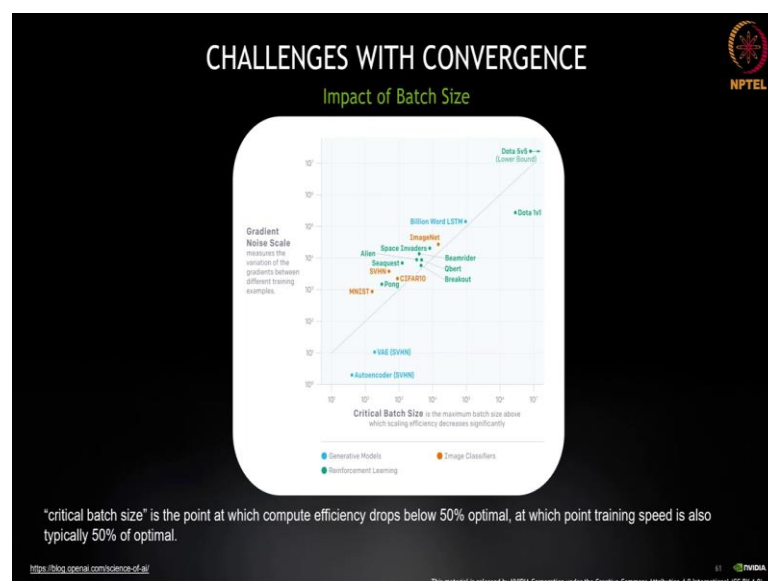
(Refer Slide Time: 00:14)



So, let us get started on the next topic we showed you how to use the distributed deep learning using two frameworks. But when you start using distributed deep learning, there are certain problems or challenges that you will face and that is what we are going to look in this particular certain slides that we have.

(Refer Slide Time: 00:38)



We see an improvement in performance when we increase the batch size that is what we showed on the first day, but it has been found that once we hit the limit which is also referred to as a critical batch size, we exhaust the possible data parallelism that we can achieve. So, you can see a paper which kind of highlights, this part if you can see there is a particular relationship that they are trying to show by increasing the batch size and it affects when you want to improve the data parallelism and it kind of exposes a particular limit which is referred to as a critical batch size.
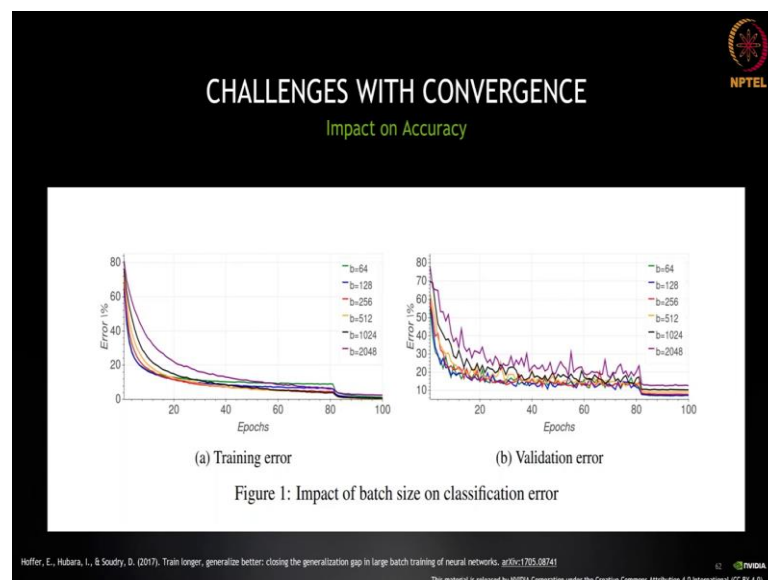
(Refer Slide Time: 01:22)

So, here you can see a critical batch size is a point at which the compute efficiency drops below 50 percent of the optimal at which point the training speed is also typically 50 percent of optimal. The critical batch size is found out to be different for different networks and that is what its being shown here and it is a function of the gradient noise scale. As you can see here your y axis is actually gradient noise scale which measures the variation of gradient between different training examples.

So, the critical batch size will be different for different networks and it is a function of gradient noise scale. So, when the noise scale is small, looking at a lot of data in parallel quickly becomes redundant whereas, if it is large we can still learn a lot from huge batch sizes and that is what the relationship is being shown in this particular paper here.

(Refer Slide Time: 02:25)



One more problem that we find with the scaling is the loss of accuracy. What you are seeing here is a graph which is part of the paper where you can see when you change the batch size actually there is a jump you can see that the accuracy is fluctuating and that is what was also observed.

So, if you change this hyper parameter, we can notice the change in the above graph right how the error actually oscillates when you increase or decrease the batch size now that is also very very critical. So, if you change anything which is also referred to as a hyper parameter, it may result into certain problems with respect to the accuracy levels.
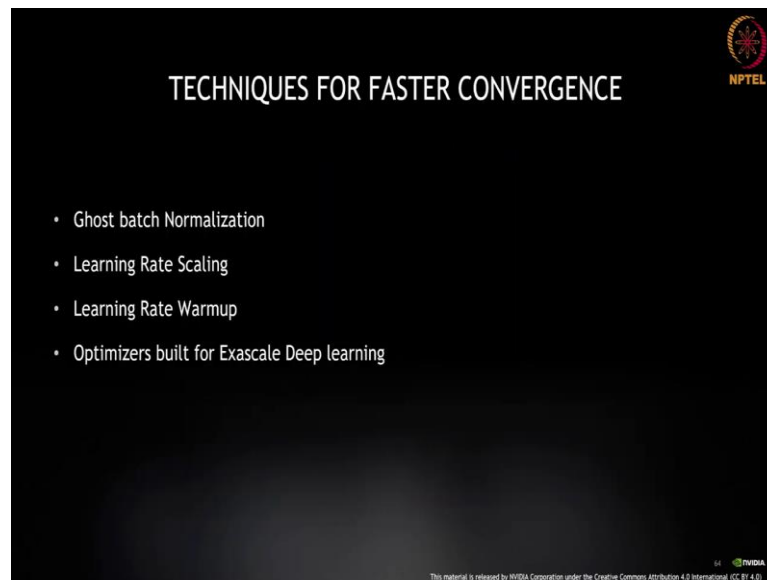
(Refer Slide Time: 03:12)



Let us look into more detail which is also highlighted in this paper where having a smaller batch size basically has a higher noise as compared to large batch size and help with getting out of local minimas which is in the case with the large batch size. Also in the above picture that you see, we see that we see two kinds of minimas sharp minima and a flat minima and these minimas are basically function of data sets when the original function is having a slight offset.

The flat minima error is primarily negligible whereas, the sharp minima function as a very large error and that is what makes it really difficult and it can create many problems and it is a function of the data set and this convergence to a sharp minima happens more frequently in the case of higher batch sizes.

And the optimizer basically can find both of them as equally good candidate. So, you can see here based on your data set this curve or this function of your loss is going to change and based on even if there was a slight offset with respect to change you did in the hyper parameter in the flat minima, it may not impact your accuracy so much.

But in case your data set was of a function in which was of type having a sharp minima even a slight change can result into huge variations and these kinds of things are what creates lot of problems and there are different ways to solve this also and we are going to look at couple of them in a bit.

(Refer Slide Time: 04:58)



So, we have seen what are the problems with scaling and the reasons behind them, but let us see some of these techniques. We are going to cover four techniques which are listed here ghost batch normalization, we are going to talk about learning rate which is another hyper parameter and usually we can change the learning rate and it is not usually constant.

There is a concept of warm up which we will discuss in more details and why it is important followed by using optimizers which are not traditionally seen in the normal deep learning world, but they become really important when we are talking about petascale level of the scaling or 2000s of GPUs as well.

(Refer Slide Time: 05:48)



So, batch normalization basically lead us add different noise to different parts of the batch independently. This increase in noise helps in learning better. This can be implemented in Horovod as its follows the MPI model. So, different GPUs have their own batches and we can train with them this improves the overall accuracy of our model.

So, you will see that if you use Horovod versus Tensorflow or the in a implementation of Pytorch if you distribute the work even though the model is the same everything is the same if the distribution strategy is different, you would realize that one of them is converging faster as compared to the other distribution strategy and this is quite common and you can reduce the impact of it by doing certain features like batch normalization.

(Refer Slide Time: 06:45)



The second example here is for faster convergence is the learning rate. You can again see the reference to the paper by scaling the learning rate basically what this paper states is that by scaling the learning rate linearly with the batch size, it has a model convergence rate and the number of epochs needed to converge would be preserved that way and our training becomes faster. The linear scaling rule it states is as follows. It states that when the batch size is multiplied by k you should also multiply the learning rate by k.

So, if you have mini batch sizes multiplied by k if you are changing your batch size in that case your learning rate should also increase by k and that is what was the observation inside this particular paper which was there.

(Refer Slide Time: 07:42)



But one side effect with the linear scaling rule we saw in the last slide is that it breaks down when the network is changing very rapidly and as you would know the rapid change of network generally happens in the early stage of training because all of your weights are random. And you are actually changing at a very higher rate and that kind of occurs in the early stage and this issue can be alleviated by a properly designed learning rate warm up a strategy of using less aggressive learning rate at the start of training and increasing it gradually.

So, you can see here for the initial first couple of epochs, the learning rate is not increased at a particular a higher way, but later on it gets increased. So, that is what its very very important and its also referred to as learning rate warm up. It has really good impact at least in the initial stages because otherwise your linear scaling that we define will be basically break down.

Finally, what we are also saying is that we can use optimizers that are built for exascale deep learning that tackle the scaling problem to a certain extent. There are different kinds of different types of optimizers which are built for exascale or petascale level deep learning like here you can see there are two graphs one of the graph is with lars which is another type of optimizer and without lars which is a traditional optimizer and you can see at the accuracy which you get is very very different right. So, and that is what it is the one that which is being observed.

So, let me again go back to the Jupyter notebook to give you.

(Refer Slide Time: 09:44)



(Refer Slide Time: 09:46)



Let me just kill this overall thing.

(Refer Slide Time: 09:52)



So, I am going to shut down all the kernels just to make sure that there is nothing which is going wrong in what I am trying to do here in this particular slide yeah.

(Refer Slide Time: 10:06)

(Refer Slide Time: 10:09)



So, at least my kernel is running now. So, we already talked about the challenges with the convergence that we saw.

(Refer Slide Time: 10:18)



And in the last era we saw.

(Refer Slide Time: 10:28)



So, we also talked about the impact of the batch sizes in the paper you would see the impact of it which. So, you would have a reference to all of these papers in general.

(Refer Slide Time: 10:29)

(Refer Slide Time: 10:36)



So, you do not need to worry about it you will get references you can download these papers and also refer to them ok its the problem again obtained ok.

(Refer Slide Time: 10:53)

(Refer Slide Time: 10:54)



So, its almost the same thing that I talked to you about.

(Refer Slide Time: 10:58)



But hopefully if Horovod runs I would be able to show you a very quick demo it does not run we inform ok.

(Refer Slide Time: 11:06)



It is running now. So, here we are using another data set this particular data set is of CIFAR which is another data set and we are basically running it and you can see here we are using horovod -np 1 which means I am practically running it at rate as in on a single GPU.

(Refer Slide Time: 11:35)



Because that single GPU is the one which we want to observe the accuracy levels ok.

(Refer Slide Time: 11:40)



(Refer Slide Time: 11:40)

(Refer Slide Time: 11:41)



(Refer Slide Time: 11:42)



So, we have got the result for the single GPU the images per second is around 42 or so, let us see if you are able to run it on the 8 GPU also and most probably it should run if it does not run here I will just stick to the theory and we can see.

(Refer Slide Time: 12:10)



So, it started running. The most important thing for here also for us is the accuracy and that is what we want to look at.

(Refer Slide Time: 12:20)

(Refer Slide Time: 12:22)



So, you can see here the accuracy for us after running on 12 epochs is 0.4434 while the accuracy here is 0.7201. So, the single GPU run gave 72 percent accuracy while the 8 GPU run gave around it gave me around only 42 percent accuracy. So, there is obviously, a loss which we describe sometime back which we can get and that is what is also highlighted that you are getting around 70 some or 80 percentage when you run it on single versus when you run it on multi GPU.

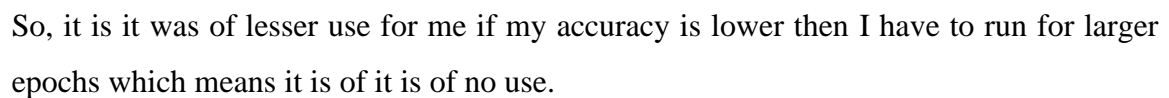(Refer Slide Time: 13:08)

(Refer Slide Time: 13:13)



(Refer Slide Time: 13:20)



So, the first technique that we mentioned for faster convergence was the use of batch normalization and for the existing code we are adding batch normalization to the same code it just that we added one more line of batch normalization here to the same code. And let us say there is a improve in the accuracy by adding this batch normalization. We are running again on 8 GPUs please see that I am running again on 8 GPUs.

So, previously by running on 8 GPUs we were able to get higher images per second with 8 GPUs I got 2,14,000 images per second versus 42,000 images per second when I was running on a single GPU, but my accuracy reduced.

(Refer Slide Time: 14:14)



So, it is it was of lesser use for me if my accuracy is lower then I have to run for larger epochs which means it is of it is of no use.

(Refer Slide Time: 14:18)

You can see here now the accuracy have increased from 40 to something in 64 by having the batch normalization. So, the by adding batch normalization my images per second decreased slightly, but I got more accuracy.

(Refer Slide Time: 14:40)



So, that is another part which needs to be seen that by adding some of these techniques you might end up reducing your overall throughput, but you what you are trying to find is the right balance of accuracy how many epochs you need right.

(Refer Slide Time: 14:52)

Then the next thing that we also talked about is changing the learning rate and you can see here we are trying to change the learning rate based on the; based on the number of GPUs here you are seeing that I am multiplying it by 0.01 and we are going to change the learning rate and not keep it make it adaptive rather than keeping it constant.

(Refer Slide Time: 15:17)



(Refer Slide Time: 15:35)

(Refer Slide Time: 15:38)



So, once we do that let us see if there is improve in accuracy or reduce in accuracy not so much. So, we can see here the accuracy increase was better than the better than the normal one. So, the learning rate did increase it, but batch normalization basically had a higher impact as compared to changing the learning rate. So, the learning rate increased from for 0.46 to 0.53 in 12 epochs and you can see here this is also really good, but still it is lesser than the learning rate.

(Refer Slide Time: 16:08)

So, we talked about the problem with learning rate which was the initial warm up which can create these kinds of problems.

(Refer Slide Time: 16:15)



To solve that you can pass parameters like Horovod dot callback and you can do a learning rate warm up and you can define that what is your warm up till how much do you want to warm up and here you can see we are saying for 3 epochs should be considered as warm ups and then later on you can do the scaling of your of the learning rate.

And let us see if it does have a impact on the accuracy. So, the normal learning rate was giving me around 53 percent accuracy.

(Refer Slide Time: 17:01)



(Refer Slide Time: 17:04)



So, you can see here again it has increased it has come back to 64 percent as compared to 46 and my images per second is also is better. So, you can see here this impact of these hyper parameters is quite a lot when you move towards distributed deep learning and if you do not do enough research and if you do not use some of these techniques.

(Refer Slide Time: 17:26)



You might end up using all the GPUs, but then you will have to use larger epochs to converge and then all of your effort in terms of reducing the time just goes away. Last, but not the least basically I am I also talked about the uses of different kinds of optimizers which are meant for these kinds of higher level of scaling like here there is a different optimizer which is called as a lamb optimizer which uses a layerwise adaptive large batch optimization to train.

And it is well known because it requires very less or little hyper parameter tuning which is very good and you can enable very large batch sizes in this without having to degrade any performance. So, rather than using the normal ones you can see here you can use this optimizer called as lamb and let us see if we are what kind of a accuracy levels we can get via using this new optimizer.
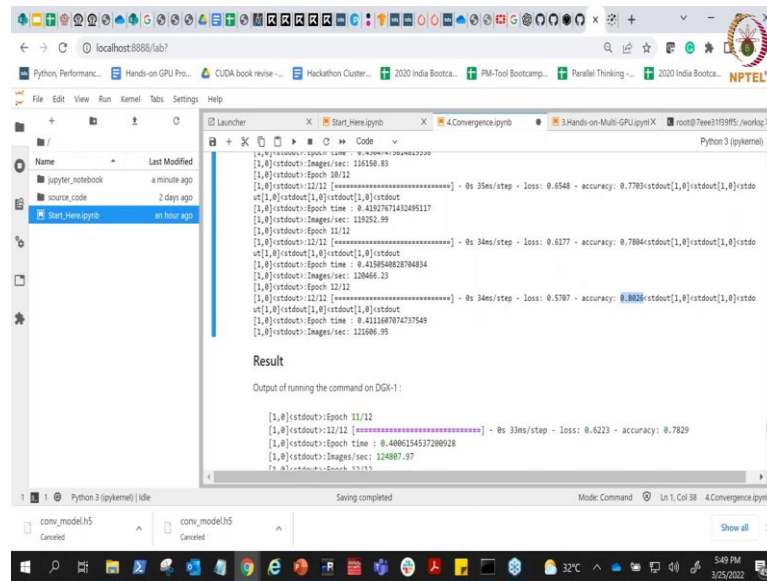
(Refer Slide Time: 18:13)



(Refer Slide Time: 18:49)

(Refer Slide Time: 18:52)



So, you can see here it is actually surpassed the capability of our traditional model also you can see here it has reached almost 80 percent accuracy which is very good it has surpass the normal optimizers and it was able to scale also not at the rate otherwise, but it was able to do many more instead of 46,000 it was able to do 1,21,000 images per second which is almost like 3 x speedup and, but also maintaining the accuracy level that you would have expected out of it.

(Refer Slide Time: 19:28)

So, if you see we started with single GPU which gave around 70 percent we use a naive implementation of multi-GPU which reduced it, but then we use techniques like batch normalization learning rate with warm up and then finally, we also looked at certain additional optimizers which are much more suited when we talk about distributed deep learning.

(Refer Slide Time: 19:50)



So, with that we are kind of done with the three series lecture that we were having on the distributed deep learning and a just a very quick recap in terms of what all things which we did, today with the increase in data set and the model sizes its almost impossible to be productive without having to use processes like GPU, but not just one processes, but many of them to make sure that we are basically working on research and not just waiting for it to finish for multiple days or weeks or months.

With that comes the challenge of different kinds of a frameworks needs to support this distribution, we saw many frameworks supporting it we saw example of running Horovod and TensorFlow Pytorch also supports it and we also saw the impact of running this frameworks on different sets of topology on different sets of hardware and what kind of a impact it can have we also demonstrated to you how the GPU computing infrastructure is not about the hardware alone.

Behind the scenes, the frameworks like TensorFlow Pytorch and all they use a highly optimized library like NCCL to do all of this effective communication in a very very low

latency high bandwidth manner. We also saw how it can choose really high speed bandwidth or a low latency interconnects like NVLink and you optimize your scaling efficiency.

Finally, we move towards even if you do all of this there are higher chances that you may end up not getting the accuracy that you are expecting out of your code for which you need to know additional techniques that we saw like batch normalization learning rate warm up and also using other kinds of optimizers which are traditionally not used when you are running it on a non distributed environment. With that we are done with the three lecture series of distributed deep learning.