

Improving the Performance of Low Bandwidth Distributed Data Parallelism

Christopher Rae

email: 110641877@ea.edin.sch.uk

Abstract

There is no doubt that networking is a bottleneck when training neural networks across multiple nodes. Having a low bandwidth or a high delay in time it takes for data to be transferred between nodes can drastically increase training time of a neural network. It has become industry standard to use fast networking with bandwidths of up to 100Gbps this makes the requirements for training across multiple nodes very expensive.

2. Research methods
3. is viable

1 Introduction

Distributed data parallelism is a widely used technique for training artificial intelligence (AI) models, especially when working with large datasets. It involves dividing the data across multiple nodes and training multiple copies of the model in parallel, which can significantly reduce the time required for training. However, the efficiency of this technique can be affected by several factors, including the bandwidth of the communication channels between the nodes.

Low WiFi bandwidths can have a significant impact on the training time of AI models using distributed data parallelism. When the communication channels between nodes have low bandwidth, it can lead to slower data transfer and synchronization between the nodes, which can result in slower training times. In addition, low WiFi bandwidths may also lead to increased latency in the communication between nodes, further diminishing the efficiency of the training process.

In this paper, we aim to explore the relationship between low WiFi bandwidths and the training time of AI models using distributed data parallelism. We will discuss the various ways in which low WiFi bandwidths can impact the communication between nodes and the overall efficiency of the training process. We will also present potential solutions for mitigating the negative effects of low WiFi bandwidths on model training time.

To provide a better understanding of the role that WiFi bandwidth plays in the training of AI models using distributed data parallelism, we will conduct experiments using a variety of datasets and model architectures. We will compare the training time of models trained using distributed data parallelism in high and low WiFi bandwidth environments to evaluate the impact of low WiFi bandwidths on the training process.

Overall, our goal is to shed light on the importance of WiFi bandwidth in the training of AI models using distributed data parallelism and to identify strategies for optimizing the training process in low bandwidth environments. We hope that the findings of this research will be useful for

practitioners seeking to improve the efficiency of their model training process and for researchers studying the optimization of distributed data parallelism for AI model training.

2 Methodology

Distributed Data Parallelism

The goal of distributed data parallelism is to split the data across the GPUs and computers in the cluster, so that the training process can be completed faster by leveraging the additional computational resources.

Here is an overview of the process:

1. First, the data and model are partitioned or "split" across the GPUs and computers in the cluster. This is typically done by dividing the data into smaller chunks and assigning each chunk to a GPU or computer in the cluster.
2. Next, the training process begins. On each GPU or computer in the cluster, the model processes the data assigned to it and computes the gradients of the loss function with respect to the model parameters.
3. The gradients computed on each GPU or computer are then averaged together, and the model parameters are updated using this average gradient. This process is known as "gradient averaging" or "gradient synchronization."
4. The process is then repeated, with each GPU or computer processing its assigned data and computing gradients, until the model has been trained.

By training the model in this way, the computation is distributed across multiple GPUs and computers, which can significantly reduce the time it takes to train the model. However, it also adds complexity to the training process, as the gradients must be averaged across the GPUs and computers and the model parameters must be kept consistent across the cluster. There are multiple ways of performing this gradient averaging depending on the architecture of your GPUs and computers. In this paper we work with a synchronous architecture (all GPUs have roughly the same compute power) for simplicity, but there is also the argument of a centralised versus decentralised topology.

In a centralized topology, the process of gradient averaging or gradient synchronization in step 3 of the distributed data parallelism process is typically performed on a central server or "parameter server." This is done by sending the gradients computed on each GPU or computer in the cluster to the central server, where they are averaged together and the model parameters are updated.

In a decentralized topology, the process of gradient averaging or gradient synchronization is performed in a decentralized manner, without the use of a central server or parameter server. This can be done using a number of different techniques, such as peer-to-peer communication or gossip protocols.

We find in this paper —[INSERT LINK](#)— that decentralized topology's are effected less by low bandwidths, and that is why we implement a decentralized architecture in our papers. To be more specific, we use our own implementation of the Baidu Ring AllReduce algorithm to average our gradients.

SGD

The training data is denoted as $\{x_i\}_{i=1}^n \subset \mathbb{R}^d$, and the model parameter vector is denoted as $w \in \mathbb{R}^d$. The goal of training is to solve the following optimization problem

$$\min_w f(w) = \frac{1}{n} \sum_{i=1}^n L(w, x_i)$$

We typically use an algorithm such as stochastic gradient descent (SGD) to solve the equation above. In which we need to calculate the unbiased stochastic gradients of f at time step t which satisfies $\mathbb{E}[g_t(w)] = \nabla f(w)$. SGD written as:

$$w_{t+1} = w_t - \eta g_t(w_t)$$

Where $g_t(w_t)$ is an estimate for the actual gradient $\nabla f(w_t)$ and η is the learning rate. For simplicity we'll let g represent $g_t(w_t)$. The convergence of SGD is largely influenced by the variance of g . In most implementations of distributed parallelism it is g which is averaged between all of the replicas.

Gradient Quantization

Gradient quantization is a technique used in distributed data parallelism to reduce the amount of data that needs to be transmitted over channels during training. It works by reducing the precision of the gradients that are transmitted between worker nodes.

A model parameter is usually stored in memory as a 32bit floating point. There are 7,800,000(38zeros) unique float32 numbers, but that said numbers like 0.000100, 0.000101 and 0.000102 are very close. We can approximate these to number to be equal, this is the core concept of quantization. When we calculate

If you have a vector $a \in \mathbb{R}^d$ where $a \sim \mathcal{N}(0.1, 0.55)$ where we know the $\max a$ and $\min a$. We can split a into equal size bins and represent each bin with either a 8bit or 16bit integer. For example, suppose we have a set of 100 bins, each representing a range of gradient values. If the gradient value for a particular parameter falls within the range of the first bin, it will be represented as a 1; if it falls within the range of the second bin, it will be represented as a 2; and so on. Then once it has been sent to the second node the node can look up what values the bin represents and replaces it with the expected value of that bin. This allows us to represent each gradient value using a single integer, rather than a floating-point value, which can significantly reduce the amount of memory required to store gradients. The difference between using 8bit or 16bit integer is the number of bins, which effect not only the accuracy of the data but how much data is sent over the network. If we use 16 bit integers there can be 65,536 bins where as with 8bit there is only 256 which is a significant reduction.

As mentioned above $a \sim \mathcal{N}(0.1, 0.55)$ this tells us that that most of the values will be located around 0.1 which means it may be more representative of a if there are more bins located around the mean than $\max a$ and $\min a$ so bins near the mean cover a lower range of values meaning they are occur at a higher frequency, bins closer to the min and the max cover a larger range of values. This type of quantization is called non-linear quantization.

Modern models can contain billions of trainable parameters, which translates to billions of gradients. If we represent all of those gradient as 256 different integers we loose a lot of accuracy so what we

do is we split gradients into chunks of around 2000 values and quantize each chunk separately this means that the first bin in the first chunk might cover a different ranges of values to the first bin in the second chunk. there are a few reasons for this, the first being that it allows us to use 8bit integers without a significant loss in accuracy, the second being that we can allocate each chunk to a specific GPU core. If one of the values in a is 7 this is a clear outlier this will effect one of the bins causing 1 bin to be wasted because of single outlier. By splitting up the data into chunks only a hand full of chunks will be effected by outliers. This type of quantization is called block-wise quantization.

Formula

Gradient Sparsification

Gradient sparsification is a technique that involves identifying and removing "zero" or nearly zero gradients from the gradients being transmitted between GPUs or computers in a distributed data parallelism setting. The goal is to reduce the size of the gradients without significantly affecting their quality, so that they can be transmitted more efficiently over the network.

In order to perform gradient sparsification we have to decide on a threshold to determine whether or not a gradient is a zero gradient or not, if the gradient is less than the threshold then it is deemed a zero gradient. In our implementation of gradient sparsification we accumulate zero gradients locally and sum them with gradients from the next iteration.

Formula

We denote sparsified gradients as $S(g, Z_t)$. We have to remember that $g \in \mathbb{R}^d$ so we let g_i be the i -th component of the vector ($g = [g_1, \dots, g_d]$). Z is a vector with the same shape as g , when the w is initialized at the start of the training process $Z_i = 0$. The function S has 2 outputs the first being the new representation of the gradient of the weights at t , v and the second being the accumulated gradients Z_{t+1}

Sparsification of the gradients is shown as follows:

$$v_{t,i} = \begin{cases} 0 & g_i + Z_{t,i} < \alpha \\ g_i + Z_{t,i} & g_i + Z_{t,i} \geq \alpha \end{cases}$$

$$Z_{t+1,i} = \begin{cases} Z_{t,i} + g_i & g_i + Z_{t,i} < \alpha \\ 0 & g_i + Z_{t,i} \geq \alpha \end{cases}$$

where α is the threshold

Larger Batch Sizes

Increasing the batch size in a distributed training setup can also improve training efficiency by reducing the frequency of the gradient averaging step. This is because each device will process a larger amount of data before the gradients are averaged and applied to the model weights. This can reduce the overhead associated with performing the gradient averaging step, and can allow the model to make more efficient use of the available computation resources.

However, it is important to carefully tune the batch size, as increasing it too much can negatively impact the convergence and accuracy of the model. This is because a larger batch size can

lead to a more noisy gradient estimate, which can make it more difficult for the model to learn effectively.

3 Results

4 Conclusion