

Chapter 14

High-Performance Computing Framework for Virtual Memory Using CNN



K. Rameshwaraiah, S. Sree Hari Raju, and K. Ashok Kumar

Introduction

Machine learning has shown impressive results in numerous image processing applications identification, natural language processing, text analytics, and safety in recent machine learning research [1–3]. With larger deep neural networks and also more training information, highly accurate models are created. When the scale of the deep learning models and also the number of training examples are multiplied, the quantity of computing needed increases proportionally [4]. It takes time to train convolutional neural networks on a computer due to the higher quantity of machine learning to compute; therefore, dispersed and great-throughput computing structures & capabilities were necessary for producing period machine techniques using great results [5]. Remote computational intelligence systems have explored requirements specifications. Operators using decentralized deeper learning technologies construct deep convolutional neural networks in a decentralized manner. Furthermore, participants must share vast supervised neural variables themselves, which ultimately generate significant communications latency, accounting significant fraction of the total time spent on dispersed machine learning [6]. As a result, the idle period of processing elements including the CPU and GPU increases, resulting in a decrease in computing resource usage [7].

K. Rameshwaraiah (✉) · S. Sree Hari Raju · K. Ashok Kumar
Nalla Narasimha Reddy Education Society's Group of Institutions, Hyderabad, Telangana, India
e-mail: ramhyd20@gmail.com

S. Sree Hari Raju
e-mail: rsv2raju@gmail.com

K. Ashok Kumar
e-mail: ashok.kondra123@gmail.com

The significant communications cost is caused by network processing and also the huge and dominant variable transmission [8]. Storage copying procedures are performed during the calculation, and the associated outgoings also are incurred when working using a multi-layer communication protocol. The number of variables in a deeper learning model grows rapidly as its size grows. For quick and huge parallel computing, ring interconnection connectivity is required. Humans offer a strategy for effectively reducing the communication cost of huge, dispersed machine learning [9].

Related Works

Our technique does this by allowing remote employees to quickly exchange supervised neural variation. The sharing process can be accomplished by creating a new mind network environment and granting access to the system storage to all employees [10]. The proposed global supervised neural main memory architecture could be used in high-performance groups coupled using high-speed networking techniques [11]. Moreover, the proposed architecture employs the RDMA technology, which avoids data communication transfer procedures across implementation and kernel-level memory buffers [12]. Finally, employing RDMA for reading/writing information from remote node storage, the proposed scheme provides a methodology for sharing supervised learning variables between dispersed employees.

Hogwild proved that machine learning activity using the same weights and separate information shards may learn DNN using asynchronous SGDs in the main storage structure without locked, which is among the globally computational intelligence optimization techniques depending on Stochastic Gradient Descent (SGD). In particular, the Dogwild architecture, which is a Hogwild modification, conducts DNN learning by asynchronously swapping weight and biases between both the masters and slave operations [13]. Whenever the master process gets varied from the slave operations, it changes the world values and distributes the revised values to every slave program at the same time. Every slave process of understanding the grades to the masters after updating a most recently obtained weight with the gradient it has discovered on its own. As a method of asynchronous SGD, the Dist Belief proposed approach Downpour SGD, which enables a huge number of modeling replicates. The Sandblaster batching optimization technique, which is a system that supports diverse networked systems, is also developed. Inside an asynchronously SGD, every DNN trained agent could acquire parameters at varying speeds without synchronization cost, maximizing the use of CPU & internet connectivity. The capabilities for a heterogeneous HPC method, which consists of CPUs & GPUs with varying requirements (clock speed, number of units, so on), could be effectively utilized to achieve optimum computing capacity.

Proposed Method

The goal of parallel processing, machine learning is to use increased technology nodes to train enormous data and huge DNN algorithms. The needed transmission power grows quickly for the size of DNN & also a dimensionality (i.e. Resolution in image information) for generating accuracy grows. As a result, low-bandwidth systems like 1Gbps Ethernet are unable to cope using the fast-growing internode supervised neural variable traffic. As a result, many businesses and research rely on high-speed network interfaces, which are commonly found in powerful computers and powerful computing platforms.

The HPC system proposed in the work comprises numerous processing elements of multiple GPUs, a different storage node with enough storage space to offer distant memory space, and also an Infiniband switching that combines those components. SMU Librarian seems to be a consumer stationary framework that offers application software for distributed software procedures and is statistically connected with the procedures during program execution. The SMU Device Controller and Infiniband Communication Layer were also kernel-level components that must be activated before the recommended process can be executed. The storage space of the nodes is provided by the SMU Servers by pooling it in preparation, or it could be provided on-demand for customer applications. The provided storage sectors of customer activities should indeed be identified in the hosting channel adaptor drivers to facilitate RDMA operations from the deeper active learning in the compute nodes. The virtual to physical address translation data for the specified web page are enrolled throughout the Signup process. Pooling storage approaches have often been favored because of the longer registration time. The SMU Device Controller seems to be a kernel-level device driver which translates a virtual address of the DNN learning cycle to the actual memory space sectors that act as the caching for the distant sharing program memory, as shown in Fig. 14.1.

Results and Discussions

The structure of the SMU Client, including its recollection pool management technique for providing the main memory buffering, can be seen in Fig. 14.2. The SMU Server is used as a user program on Linux. The multi-threaded design of the SMU Server allows it to spin to process requests from numerous SMU Peripheral Devices at the same time. The messaging reception threading temporarily stores the response message being sent by SMU Peripheral Devices in the messaging rings. The messaging recipient process wakes all the demand processes in the memory pool followed by information collection in the rings. Perhaps one of those, in a circular queuing pattern, accepts one request message from the messaging rings, analyzes it, and would then return to the ready queue whenever the messaging process is complete. The SMU Server uses the buddy memory control method, which would be

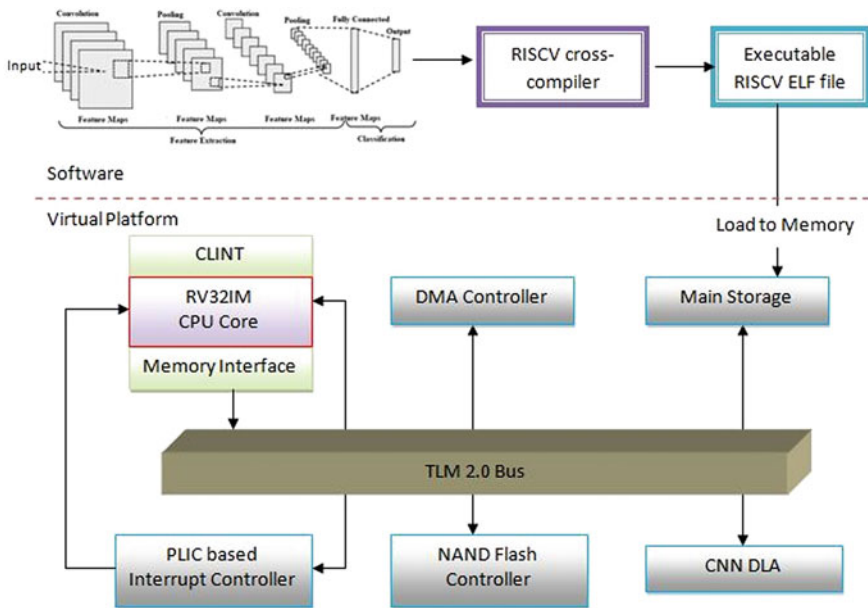


Fig. 14.1 System architecture

typically used for normal pages assignment in a Linux kernel, to handle given information as a pool. With less than 4 KB & even a maximum level of 32 GB, a storage block length could be customized, a maximum length can also be changed. Libibverbs have been used to build engaging lessons. The SMU provides the network with the remote management keys needed to the SMU system Driver for readily accessing main storage buffers during an address space creation stage.

The SMU has been used to allocate and deallocate memory space across distributed machine learning employees, as seen in Fig. 14.3. Firstly, the master’s deep learning person uses the SMU API to construct a main memory buffer on the SMU Controller. The expert task distributes SHMKey for the storage space for the workers establishing it. Additional workers which would like to be allotted a main storage buffer produced by an expert employee need to have their channel of communication. The person who gets the SHMKey from the supervisory employee transmits the SHMKey including a distributed memory allocation application. If the employee asks for the same sized sharing memory data allocation and delivers the very same SHMKey, the SMU Server uses the Encryption Key for the sharing buffer cache. After the assignment operation has finished, the SMU User’s shared memory buffer could be distributed across the dispersed deep learning workers.

Table 14.1 displays the variable sizes and computation times of six CNN models that are assessed. Because batch length represents several images to be learned, changing the batch size affects the computational effort. Caffe needs substantial memory resources as the batch size grows, hence there seems to be a maximum

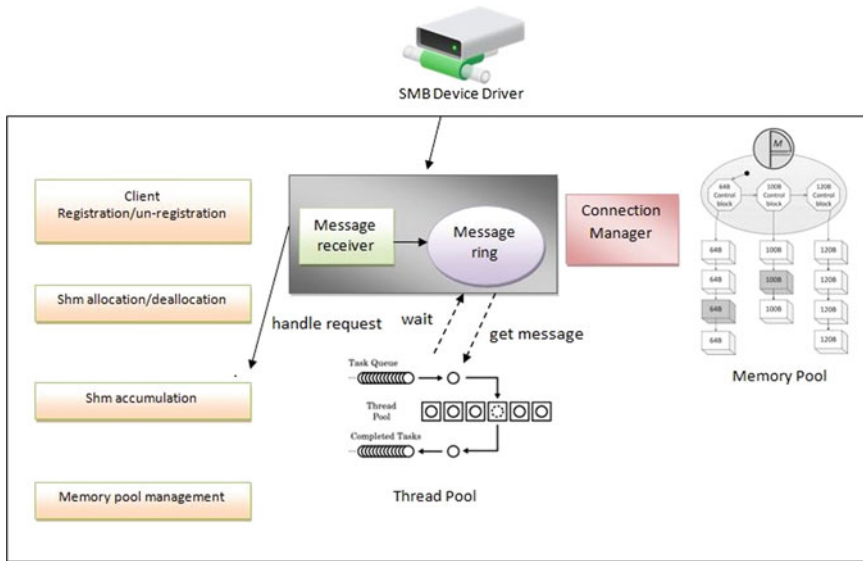


Fig. 14.2 Architecture of the SMU server

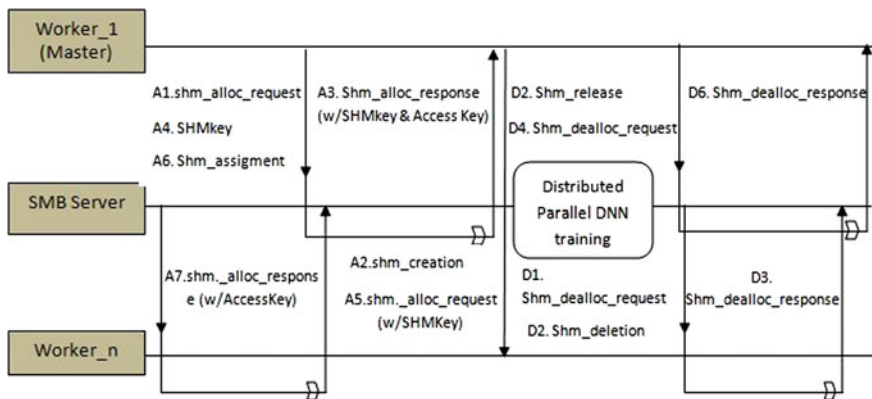


Fig. 14.3 Observational research on sharing data storage

batch length at which Caffe could train every model in an execution. In terms of variable length and computational cost, the proposed CNN models have a wide range of values. The time it would take to compute a variable is never dependent on the magnitude of the variable. The VGG16 contains the most parameters (140 million), yet it takes less time to compute than the other four methods.

In the parallel situation, the connection times of MPI and SMU are compared in Figs. 14.4 and 14.5. During every iterative process of DNN coaching, the sequences situation does not consider data processing overfitting. As demonstrated in Fig. 14.4,

Table 14.1 ML performance measures

| Network specification | CNN | | | | | |
|-----------------------|--------------|--------------|-----------|---------------------|------------|--------|
| | Inception_v1 | Inception_v3 | Resnet_50 | Inception_resnet_v2 | Resnet_152 | VGG16 |
| Variable (million) | 13.5 | 24.6 | 31.9 | 57.2 | 66.6 | 139.1 |
| Variable size (Mb) | 51.21 | 94.32 | 121.14 | 214.61 | 253.64 | 527.26 |
| Compute time (ms) | 190 | 451 | 234 | 295 | 219 | 191 |
| Batch size | 51 | 33 | 19 | 7 | 8 | 65 |

as the model size & number of staff increased, an MPI product’s transmission rate increased rapidly, whereas the SMU product’s transmission rate gradually increases. Figure 14.4 depicts general trends, while Fig. 14.5 compares the MPI and SMU techniques by adjusting the amount of machine learning workers for each of the six deep learning techniques. In all scenarios, it has been shown that the SMU technique is 1.1 to 4.2 times faster than the MPI approach.

For any dataset and any number of staff, the SMU work demonstrated the MPI approach. The SMU technique has a transmission time of 0 ms up to 24 employees in an Inception v1 using a very tiny design capacity, & 45 ms, even the number of staff is increased to 32, which would be 7 times faster than the MPI approach. In the instance of Inception v3, which has a shorter variable size and a longer computational cost, the calculation time hides the entire calculation time of the SMU technique. However, whenever the number of processes is increased to 32, the MPI product’s transmission rate exceeds the computational effort by 2.3 times. The variation in

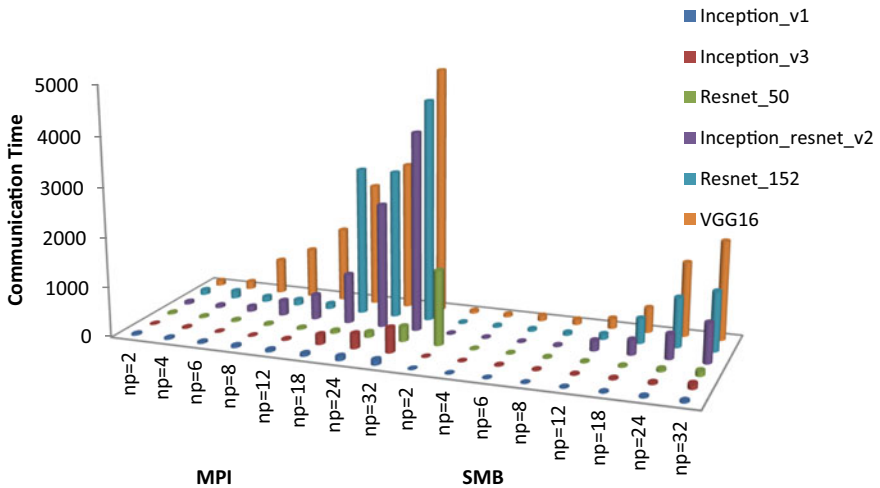


Fig. 14.4 Comparison of MPI-based and SMU

Fig. 14.5 Comparison of transmission times for each type

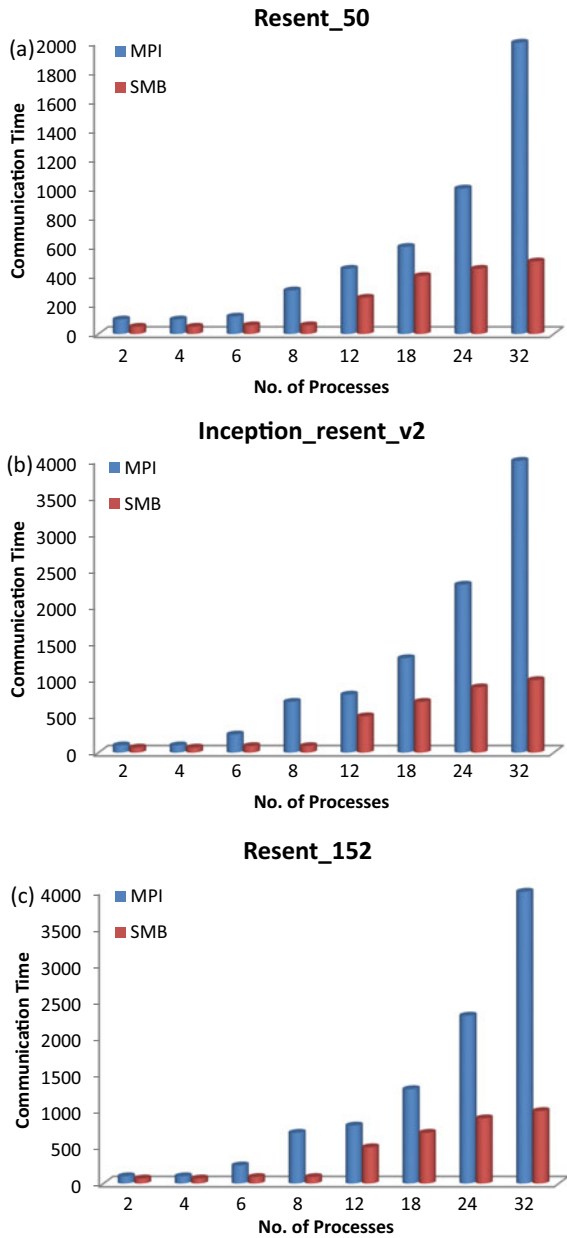
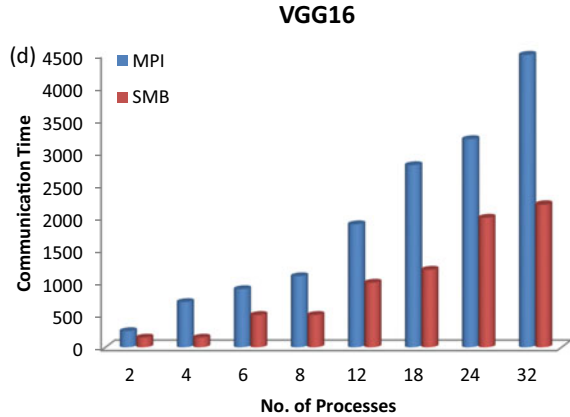


Fig. 14.5 (continued)



the true communication gap between the MPI & SMU system reduces as a model size rises, but in the particular instance of Resnet 50, 4.8% stronger in an Original conception Resnet v2 framework, & 2.0 times higher in a Resnet 152 & VGG16 prototype, the SMU technique is 5.2 times faster than the MPI technique. The SMU-based machine learning system's average of 1 repetition number of training has increased by 3.5 times.

Conclusion

The proposed SMU makes use of media transmission memory to transport variables from the memory of the local computer to the memory of the remotely distributed node, and interchange DNN variables by learning to write the remote distributed memory buffer. Storage duplication and internetwork execution have a significant influence on communication latency. The proposed SMU-based system was 2.1 times quicker than the MPI-based system in a sequential situation for asynchronous SGD. The communications with the SMU-based dispersed DNN training platform are two to seven times more efficient than using an MPI-based system in a concurrent scenario of asynchronous SGD, where computational resources overlapped. Only a minimal predictive performance is validated because the results of these experiments are achieved using a single SMU Server. The amount of SMU Server, on the other hand, could be raised to boost efficiency, and also capacity could be calculated according to the capacity of models & also the number of supervised neural employees. The suggested SMU-based method lays the groundwork for expediting the networked analysis of information DNN learning in a powerful computational environment along with high speed and accuracy.

References

1. Zhang, X., Zhang, T., Lu, J., Fu, X., Reveriano, F.: The effect of high-performance computer on deep neural network. *Eng. Sci.* **15**, 67–79 (2021)
2. Singh, A., Prakash, S., Kumar, A., Kumar, D.: A proficient approach for face detection and recognition using machine learning and high-performance computing. *Concurrency and Comput.: Practice Exper.* **34**(3), e6582 (2022)
3. Haseeb, M., Saeed, F.: High-performance computing framework for tera-scale database search of mass spectrometry data. *Nature Computat. Sci.* **1**(8), 550–561 (2021)
4. Jin, S., Li, G., Song, S. L., Tao, D.: A novel memory-efficient deep learning training framework via error-bounded lossy compression. In: *Proceedings of the 26th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, February, pp. 485–487. (2021)
5. Lima, A.L.D.C.D., Aranha, V.M., Carvalho, C.J.D.L., Nascimento, E.G.S.: Smart predictive maintenance for high-performance computing systems: a literature review. *J. Supercomput.* **77**(11), 13494–13513 (2021)
6. Deng, C., Sui, Y., Liao, S., Qian, X., Yuan, B.: GoSPA: an energy-efficient high-performance globally optimized sparse convolutional neural network accelerator. In: *2021 ACM/IEEE 48th Annual International Symposium on Computer Architecture (CA)*, June, pp. 1110–1123. IEEE (2021)
7. Pandey, S., Nagwani, N.K., Verma, S.: Aspects of programming for the implementation of convolutional neural networks on multisystem HPC architectures. *J. Phys. : Conf. Series* **2062**(1), 012016 (2021). IOP Publishing
8. Patel, S., Liu, T., Guan, H.: free lunch: compression-based GPU memory management for convolutional neural networks. In: *2021 IEEE/ACM Workshop on Memory Centric High-Performance Computing (MCHPC)*, November, pp. 1–8. IEEE (2021)
9. Latchoumi, T.P., Parthiban, L.: Quasi oppositional dragonfly algorithm for load balancing in cloud computing environment. *Wireless Personal Commun.* 1–18 (2021)
10. Yao, C., Liu, W., Tang, W., Hu, S.: EA: energy-aware adaptive scheduling for CNN inference on high-performance GPUs. *Future Generation Comput. Syst.* (2022)
11. Balamurugan, K.: Metrological changes in surface profile, chip, and temperature on end milling of M2HSS die steel. *Int. J. Mach. Mach. Mater.* **22**(6), 443–453 (2020)
12. Chang, S.E., Li, Y., Sun, M., Shi, R., So, H.K.H., Qian, X., ..., Lin, X.: Mix and match: A novel FPGA-centric deep neural network quantization framework. In: *2021 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*, February, pp. 208–220. IEEE (2021)
13. More, N., Galphade, M., Nikam, V. B., Banerjee, B.: High-performance computing: a deep learning perspective. In: *Deep Learning and Edge Computing Solutions for High-Performance Computing*, pp. 247–268. Springer, Cham (2021)