CFD Simulation and Optimization of the Cooling of Open Compute Machine Learning "Big Sur" Server

Mangesh Dhadve, Jimil M. Shah, Dereje Agonafer The University of Texas at Arlington Box 19018, 500 W. First Street RM. 211, Woolf Hall Arlington, Texas, USA

Email: mangesh.dhadve@mavs.uta.edu

ABSTRACT

In recent years, there have been phenomenal increases in Artificial Intelligence and Machine Learning that require data collection, mining and using data sets to teach computers certain things to learn, analyze image and speech recognition. Machine Learning tasks require a lot of computing power to carry out numerous calculations. Therefore, most servers are powered by Graphics Processing Units (GPUs) instead of traditional CPUs. GPUs provide more computational throughput per dollar spent than traditional CPUs. Open Compute Servers forum has introduced the state-of-the-art machine learning servers "Big Sur" recently. Big Sur unit consists of 4OU (OpenU) chassis housing eight NVidia Tesla M40 GPUs and two CPUs along with SSD storage and hotswappable fans at the rear. Management of the airflow is a critical requirement in the implementation of air cooling for rack mount servers to ensure that all components, especially critical devices such as CPUs and GPUs, receive adequate flow as per requirement. In addition, component locations within the chassis play a vital role in the passage of airflow and affect the overall system resistance. In this paper, sizeable improvement in chassis ducting is targeted to counteract effects of air diffusion at the rear of air flow duct in "Big Sur" Open Compute machine learning server wherein GPUs are located directly downstream from CPUs. A CFD simulation of the detailed server model is performed with the objective of understanding the effect of air flow bypass on GPU die temperatures and fan power consumption. The cumulative effect was studied by simulations to see improvements in fan power consumption by the server. The reduction in acoustics noise levels caused by server fans is also discussed.

KEY WORDS: Open Computer Server, Machine Learning, Big Sur, Hardware, CFD Simulation, Air flow optimization

INTRODUCTION

In 2015, at Open Compute Project (OCP) summit Facebook unveiled its first GPU server Big Sur, a purpose-built hardware for machine learning applications. Facebook have been developing Artificial Intelligence software for a considerable time with off the shelf equipment, but Facebook realized that to tackle these problems at large scale, they need to develop their own hardware and this need gave birth to the

first GPU-based server to train the neural network. Its current design supports up to eight PCI-e Gen3 full heights, dual slot cards with each card hosting up to 300W. It is designed to support standard SSI-EEB compliant motherboard and eight 2.5-inch SATA drives. This server is Open Rack V2 compatible server to meet OCP standards.

One of the most important things about this server is that its design is modular. Each component such as a motherboard, FPGA board and SSD backplane is connected to each other using cables and can be moved around as per the application. It is designed around NVIDIA Tesla M40 GPU, but it also supports Intel Xeon Phi, AMD FirePro cards. Facebook uses NVIDIA Tesla M40 cards with TDP of 250W. Therefore, we have used the same card in this study. Big Sur uses two Intel Xeon E5 CPUs with TDP of 135W. It is easily serviceable by adopting swappable fans, thumb screws, removable trays.

Entire power footprint of Big Sur is around 2.5kW, which makes it one of the highest power density servers. This highest power density also necessitates in improvised thermal solution to reduce cooling costs and improve data center efficiency at the same time.

According to the report of Berkeley Lab, in 2014, data centers in the USA alone consumed 70 billion kWh of energy that is nearly 1.8% total electricity consumption in the United States. This report also shows electricity consumption in data centers increased by 4% from 2010 to 2014. This report projects energy consumption in 2020 will be nearly 73 billion kWh [1]. In the same year, a survey conducted by Uptime Institute shows data center's average PUE across the globe is 1.7 which means data centers are spending 0.7 W extra per 1 W of computing power and this extra energy goes mostly to cooling equipment (energy loss) [2]. Improving cooling efficiency will help in huge savings in cooling costs. In the typical air-cooled server, cold air is passed over electronic components which dissipate heat to the air. This heated air is then extracted through the server fans and it exits towards the hot aisle. Improving air flow conditions, the surface area available for heat transfer and proper ducting can reduce the air flow requirements hence power savings in cooling equipment.

MODELLING

Initially, much of the data on a physical model of the server was not available, at that time we used Autodesk Fusion 360 to reverse engineer photographs of the server using canvas drawing feature as shown in figure 1. We used SolidWorks to create a 3D model of the server.

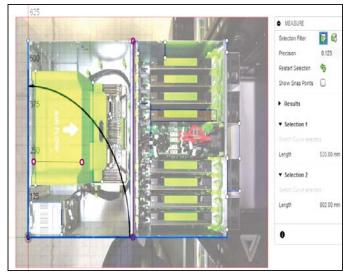


Figure 1: Extracting dimensions from image using Autodesk Fusion 360

When the 3D model of the server became available on OCP wiki, we compared it with our model and we found that our model was accurate up to \pm 2 mm.

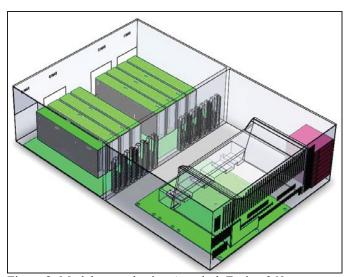


Figure 2: Model created using Autodesk Fusion 360

An original model which was obtained from OCP contained all the details needed except third party components like processors, DIMMs and PCIe cards. We used data available from Intel and OCP to model processor and PCIe cards as shown in figure 3. ANSYS SpaceClaim was used to reverse engineer .STL files.

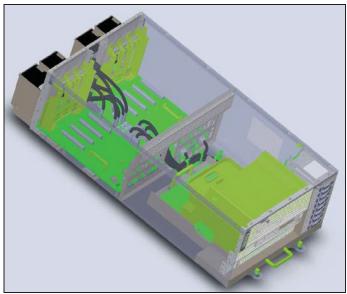


Figure 3: Model obtained from OCP

As shown in figure 4, we created a new model in SolidWorks from scratch with reference to the model shown in figure 3. We focused solely on major components creating an obstruction to airflow or modifying direction. The velocity of air depends on Air Duct, Chassis, GPU guide etc. This was an important step because it reduced the size of our model and ultimately, computational time.

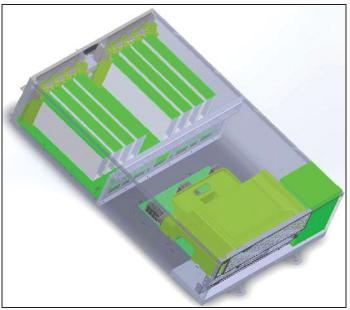


Figure 4: Model developed from the reference of OCP model

The final model is created in SolidWorks and imported in 6SigmaET as .STL file. In the model shown in figure 4, we did not model electronic components such as CPU, Fans, Heatsinks and GPU chips in SolidWorks. Those components were modeled directly in 6SigmaET as shown in figure 5.

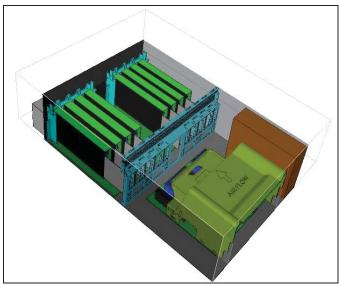


Figure 5: Model used for thermal simulation

In the thermal model, we focused on the major heat dissipating components such as DIMMs, Processors, and GPUs to create a compact model and saving a valuable computational time.

Figure 6 shows an actual model of Intel Xeon E5 package. This model cannot be used directly in the thermal model as it contains all details which do not affect thermal performance significantly. This model is simplified and made compact as shown in figure 7 in 6SigmaET, which contains all necessary components to replicate similar performance. Each package consumes 135 W power.

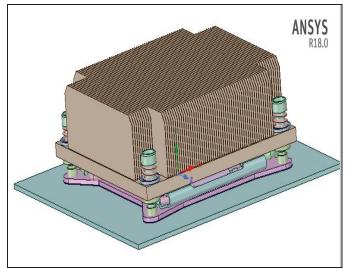


Figure 6: Intel Xeon E5 model

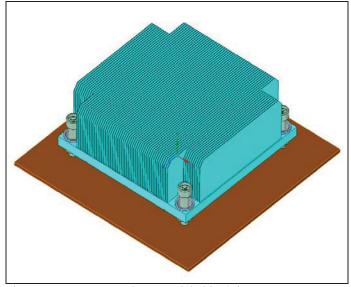


Figure 7: Processor package modeled in 6SigmaET

Big Sur uses total 8 fans of 92mm*92mm*38mm. Two fans are arranged in series so it creates total four fan modules, each module consisting of two fans in series. When trying to model them separately giving fan curve, solver struggles to settle on an operating point on fan curve for each fan. This is due to constantly changing conditions at their inlet and cells between them. In our study, we increased the cell count between two fans by changing the distance between the two fans, however the results were same. As two fans are in series, we can combine their fan curves to model them as one fan and it will still give the same performance. Theoretically, we should get double pressure difference, although in practice it is not possible due to angular component of airflow at the exhaust of fan. However, this can be minimized by guiding angular component of airflow back in the air stream by introducing guide vanes [4]. Due to the lack of knowledge of fans used, after iterations, we selected Delta Electronics' AFB0912UHE-A fan with a maximum air flow of 160.22 CFM [6]. For the selection of fan, required airflow was established by the following formula

Airflow in CFM= 154*Heat Load(kW)*(Reference Pressure/Local Pressure)

Figure 8 shows the theoretical fan curve for the fans in series and for a single fan.

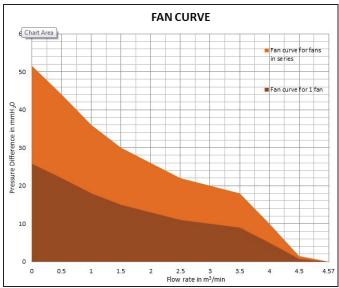


Figure 8: Fan curve for fans in series

METHODOLOGY

A CFD model is generated and baseline simulation is carried out to establish thermal parameters such as die temperatures, air velocity, hot aisle temperature etc.

In existing CFD model, air is flowing over CPUs and DIMMs then it is going towards GPUs. Therefore, thermal shadowing is occurring. Changes were made in the air duct to reduce thermal shadowing by bypassing the air.

Baseline Simulation

For this study, we used air inlet temperature of 20^oC as per ASHRAE guidelines [3, 5] and used Standard K-Epsilon turbulent model available in 6SigmaET.

In the first test case, we used inlet air temperature of 20^{0} C and fixed fan speed of 6000 RPM. Thermal properties are assigned appropriately as per the material components data available. Figure 9 shows the top view of the server with the temperature distribution of various components. In this figure, a cover of PCIe cards, their supports are hidden while the transparency of the air duct is changed.

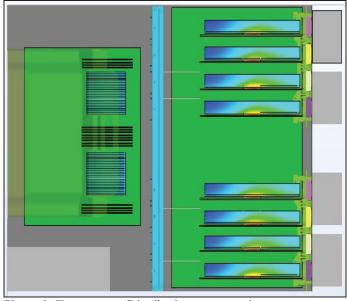


Figure 9: Temperature Distribution across various components

Figure 10 shows streamline plot obtained in baseline simulation from the top, while figure 11 shows this plot from a side view. As observed, much air is passing through space between GPUs, gap available for cable and from the top of the components.

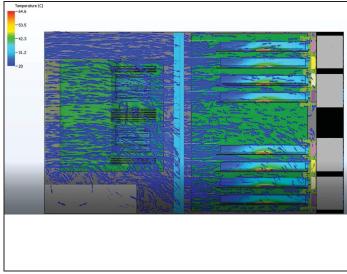


Figure 10: Streamline plot (Top View)

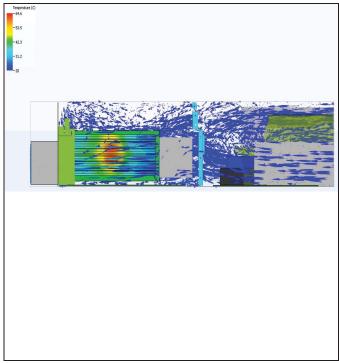


Figure 11: Streamline plot (Side View)

Grid Independence

In this study, temperatures of all the GPUs are taken into the account for grid independence. Because we observed that even if some GPU temperatures are constant during the mesh sensitivity analysis, some GPU temperatures are changing. We achieved grid independence at nearly 82 million cell count from thereafter we observed constant temperature across all GPUs. Cell count began from 35 million as recommended by Future Facilities Inc.

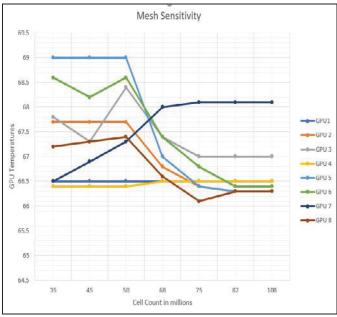


Figure 12: Mesh Sensitivity

Modifications

As observed from figure 10 and 11, middle partition creates significant resistance to the airflow. Also, Air flowing through the air duct carries on heat from CPUs and DIMMs to the GPUs, hence reducing effective convective heat transfer from GPU heatsink to the air. This study focuses on minimizing airflow resistance created by middle partition and reducing thermal shadowing effect by modifying air duct design to create a passage for air bypass.

Allowing adequate airflow to all critical components is an important factor that needs to be considered while making changes in the system. Excessive bypass air will increase the temperature of CPUs and DIMMs which is undesirable. Optimum air flow must be provided to the CPUs since they have a considerable amount of heat energy generated during operation.

Air flowing through DIMMs and part of the air flowing over CPUs is directed over the less significant components of FPGA board. While part of air over CPU and bypass air is mixed and utilized in GPUs for heat dissipation.

Figure 13 and Figure 14 shows existing air duct while figure 15 and 16 shows modified air duct.

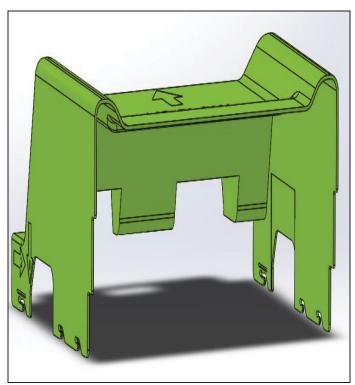


Figure 13: Existing air duct

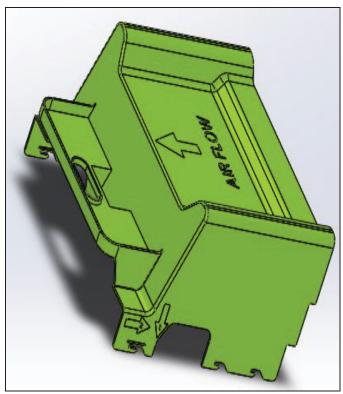


Figure 14: Existing air duct

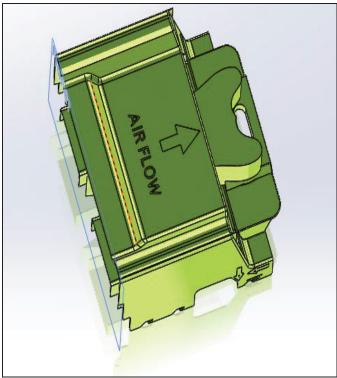


Figure 15: Modified air duct

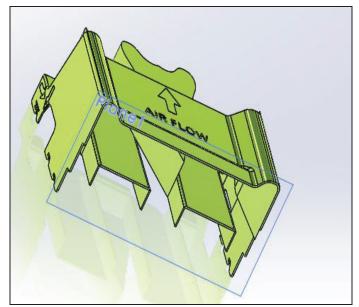


Figure 16: Modified air duct

As shown in figures, in the modified air duct, air bypassed and the middle partition in the air duct and walls are designed to reduce area available for the air flow which in turn helps to impart venturi effect on to the air flow. Also, modifications in walls are made in such a way that walls impart direction momentum to air flow. Figure 17 shows the results for modified partition with temperature distribution.

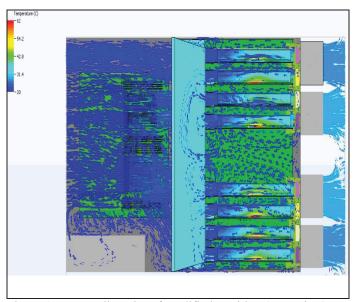


Figure 17: Streamline plot of modified partition (Top View)

Furthermore, to reduce the air flow resistance presented by middle partition, it is modified. Figure 18 shows existing partition while figure 19 shows modified partition.

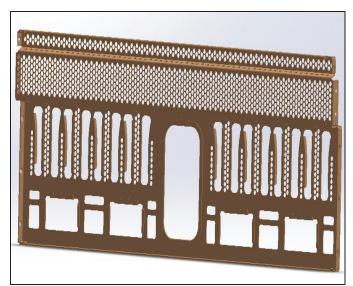


Figure 18: Existing partition

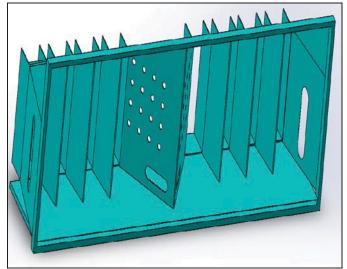


Figure 19: Modified partition

Modified partition also eliminates the need for attachments required for PCIe cards to hold them in position as it provides a guide on which card can rest. This partition instead of just giving passage for the airflow guides air into each card while presenting much-reduced resistance. Holes and slots are provided for bypass air required to cool non-significant components. Provisions have been made through slots for space needed for cables. To minimize the vortex formation in trailing edges of the holes and slots, fillets and/or chamfers are provided

RESULTS

After careful study of airflow patterns, changes were made in the air duct and partition to reduce resistance and thermal shadowing effect. This section compares the results of this modification with the baseline simulation results. We observed reduced temperature in both cases, although temperature reduction in case of air duct modifications was

not significant. During the CFD analysis, all CPUs and GPUs were clocked at 100% power levels because it will easily allow components to reach a critical temperature.

Figure 20 represents the temperature difference observed in baseline simulation and modified server.

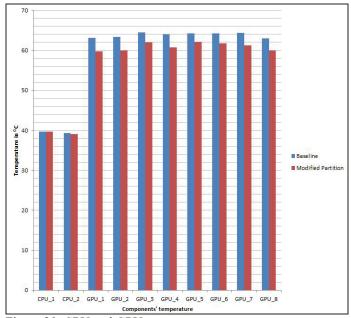


Figure 20: CPU and GPU temperatures

Figure 21 represents the fan power consumption while figure 22 shows fans flowrates in baseline simulation and modified server.

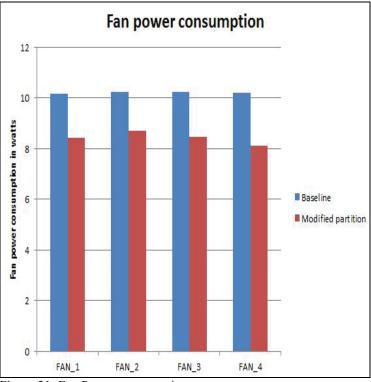


Figure 21: Fan Power consumption

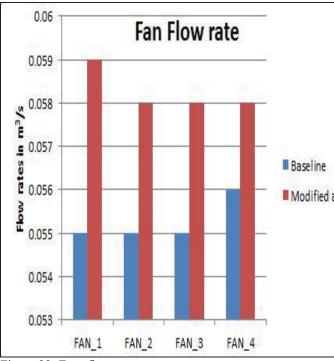


Figure 22: Fans flow rates

Due to reduction in static pressure, we observed reduction in fan noise as shown in figure 23. Relation between static pressure and fan noise is as follows [7]

$$L_n = 67 + 10 \log(s) + 10 \log(p)$$
 Where,

L_n= sound power level in dB

S = rated motor power in kW

P= static pressur in N/m²

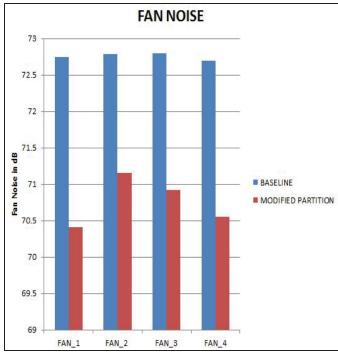


Figure 23: Fan noise

CONCLUSION

The cooling performance of Open Compute Big Sur server was evaluated through CFD analysis and 20°C inlet air temperature. The resulting data from baseline simulation was then used to modify components of the server to reduce thermal shadowing and flow resistance.

The savings achieved through these changes are as follows,

- 1. Fan power reduction- 15% to 20.8%
- 2. Temperature reduction- 3.5°C to 1.2°C

This study can be inferred at data center level to determine total savings achieved in cooling costs.

FUTURE WORK

This study focuses on CFD study and results of Big Sur server. In future, experimental testing should be done to verify data obtained through CFD analysis. Furthermore, this study is concentrated on server fitted with NVIDIA Tesla M40 cards and Intel Broadwell processors, being a modular server, Big Sur can support different PCIe cards as well as Haswell processors. Further analysis can be done with different PCIe cards such as Intel Xeon Phi, AMD FirePro etc.

Today, thermal engineers are giving more attention to oil cooling. Big Sur, being one of the heavyweight champions, in terms of power footprint, CFD analysis of this server when subjected to oil cooling can be done. With the help of oil cooling, even more powerful processors can be used to improve the computational capacity of Big Sur, which will even require using different Rack as current rack imposes power limits.

Most of NVIDIA's GPU accelerator cards are provided with passive cooling utilizing heatsink to facilitates heat transfer from GPU to the air. By utilizing active cooling, we can reduce the size of the heatsink and ultimately the size of the PCIe card and server. Thus, it can save valuable real estate at the data center level or it will be possible to add more GPUs in the same available space.

ACKNOWLEDGMENTS

The authors would like to acknowledge William Tsu (NVIDIA), Alan Chang (Quanta), Mariz Lea (Intel) for their assistance in providing needed data on various components.

The authors would like to thank John Stuewe (Open Compute Project) for providing important information on detailed CAD model of the server. Also, authors are greatly indebted of Matt (Future Facilities) for his important contribution in helping in various problems with CFD model

REFERENCES

 Arman Shehabi et. al. "United States Data Center Energy Usage Report", Ernest Orlando Lawrence Berkeley National Laboratory, U.S. Department of Energy, June 2016 [Online] Available at

https://eta.lbl.gov/publications/united-states-data-center-energy-usag

2. Uptime Institute, "2014 Data Center Industry Survey." [Online]

Available at

https://journal.uptimeinstitute.com/2014-data-center-industry-survey

- 3. ASHRAE Datacom Series 3, The green grid, White paper 18, pp. 14, 2010
- 4. Gareth Jones, "Using fans in series and parallel: performance guidelines", ebm-papst Automotive and Drives (UK) Ltd.
- 5. Jay Park, "Data Center V1.0", Open Compute Project [Online]

Available at

http://www.opencompute.org/assets/Uploads/DataCenter-Mechanical-Specifications.

6. AFB0912UHE-A Fan Specifications, Delta Electronics Inc, July 17, 2007 [Online]

Available at

http://www.deltafan.com/Download/Spec/AFB0912UHE-A.pdf

- 7. Fan Noise Power Generated, Engineering Toolbox, http://www.engineeringtoolbox.com/fan-noise-d_61.html
- 8. TF06-HeatLoad_CFM, Innovative Research, Inc [Online] Available at

 $http://titleflow.com/assets/files/titleflow/TF06Heatload_cf\\ m.pdf$