
Enhancing Data Center Sustainability with a 3D CNN-Based CFD Surrogate Model

Soumyendu Sarkar^{†*}, Avisek Naug[†], Zachariah Carmichael[†], Vineet Gundecha[†],
Antonio Guillen[†], Ashwin Ramesh Babu, Ricardo Luna Gutierrez

Hewlett Packard Enterprise (Hewlett Packard Labs)

soumyendu.sarkar, avisek.naug, zachariah.carmichael,
vineet.gundecha, antonio.guillen, ashwin.ramesh-babu, rluna
@hpe.com

Abstract

Thermal Computational Fluid Dynamics (CFD) models analyze airflow and heat distribution in data centers, but their complex computations hinder efficient energy-saving optimizations for sustainability. We introduce a new method to acquire data and model 3D Convolutional Neural Network (CNN) based surrogates for CFDs, which predict a data center’s temperature distribution based on server workload, HVAC airflow rate, and temperature set points. The surrogate model’s predictions are highly accurate, with a mean absolute error of 0.31°C compared to CFD-based ground truth temperatures. The surrogate model is 3 orders of magnitude faster than CFDs in generating the temperature maps for similar-sized data centers, enabling real-time applications. It helps to quickly identify and reduce temperature hot spots(7.7%) by redistributing workloads and saving cooling energy(2.5%). It also aids in optimizing server placement during installation, preventing issues, and increasing equipment lifespan. These optimizations boost sustainability by reducing energy use, improving server performance, and lowering environmental impact.

1 Introduction

Computational Fluid Dynamics (CFDs) are commonly used to simulate data center (DC) thermo-fluid dynamics for energy-efficient design. A significant challenge for CFD tools lies in building and then simulating models of data centers, which can house thousands of computing equipment such as servers, network devices, racks, sensors, and intricate cooling systems. However, the time-consuming nature of CFD for large data centers is often incongruent with the real-time demands of modern data centers, where timely decision-making can lead to substantial energy savings, improved performance, and reduced environmental impact.

A CFD surrogate model, that can provide results in real time, can optimize workload distribution to minimize data center hotspot temperatures Ilager et al. (2020); Lin et al. (2022); Jin et al. (2023); Fresca and Manzoni (2022). By integrating this model into current data center control software, designers can enhance control algorithms for workload scheduling and cooling optimizations, resulting in greater cooling efficiency leading to diminished carbon footprint. However, most of these approaches have made extremely simplifying assumptions for the premise of the problem: Chen et al. (2021) assumes that data center cabinets can be randomly simulated by placing rectangular grids on a

*Corresponding author

†These authors contributed equally

2D map. Furthermore, these surrogate models only generate a 2D temperature output map that has little consequence for hotspot reduction purposes.

We need an approach that can process real-world CFD data from custom data center configurations (Figure 1 left), and generate surrogate models that can "accurately" predict the temperature at any point in space of the 3D volume (Figure 1 middle). This "effectively allows" downstream applications (Figure 1 right) like temperature hotspot reduction for a given workload using an iterative optimization approach.

2 Proposed Solution

This paper applies a 3D Convolutional Neural Network (CNN) model for predicting data center 3D temperature distribution that can lead to faster inference suitable for real-time applications. The process of mapping raw CFD data to inputs and labels for training the 3D CNNs is a complex process. We developed a detailed methodology for acquiring and processing this data to create 3D input and output images for the surrogate model, as discussed next. An overview of this approach is illustrated in Figure 2.

CFD Data acquisition: We collected CFD data of the data center model with $M \times N$ arrangement of IT Cabinets (Figure 2 left) with each cabinet having a custom arrangement of an HPE Flex Fabric 5920 Switch, 4U HPE Nimble Storage HF20 and 12 HPE DL380 servers. Similar to real data centers, we considered the absence and presence of cold-aisle containment. The data acquisition process was conducted by the commercial CFD software, 6SigmaDCX 6SigmaDCX | Future Facilities (2023) automated using *6Sigma CommanderTM*: a CLI was run using a Python script to execute thousands of simulations in batches.

3D Input Data: We mapped the raw data to voxel tensors, to generate an input representation, which captures the data center's geometry and component details with each channel indicating one physical property. In our application, it includes the power of IT equipment, the ACU set point, and the ACU fan speed. The voxel tensor represents the exact 3D position and structure of components like cabinets and ACUs as specified in 6SigmaDCX. These components are depicted in a quantized room size of $64 \times 48 \times 64$ voxels from an original size of $6 \times 4.5 \times 6$ meters. This quantization ensures information accuracy and compatibility with the 3D CNN's dimensions. The metrics are volume-normalized, indicating the dispersion of each metric over its volume. Specifically, the channels (Figure 2 center) include:

- **Power:** Represents the volume-normalized power in kW of IT components, with non-component areas set to 0 kW.
- **ACU Set Point:** Encodes the normalized and negated ACU temperature in Celsius for ACU supplies and slotted grilles, with non-component areas set to 0.
- **ACU Fan Speed:** Details the volume-normalized fan speed as a percentage for ACUs and slotted grilles, with other areas set to 0

These channels are concatenated to create each input tensor $x \in \mathbb{R}^{C \times D \times H \times W}$ where C is the number of channels (three in our case), D is the room depth (64 in our case), H is the room height (48 in our case), and W is the room width (64 in our case). We denote a set of input tensors as, $X = \{x_1, \dots, x_N\}$ where N is the number of samples.

3D Output Data: The surrogate model aims to predict 3D temperature heatmaps (Figure 2 right) for each voxel in a data center based on input data. Here, we describe the construction of these

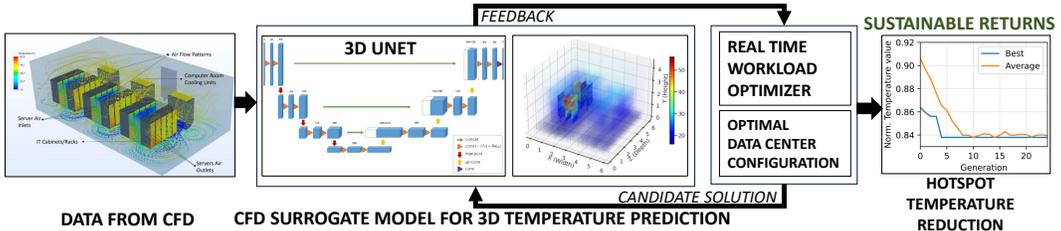


Figure 1: Enhancing Sustainable Data Centers using a Predict and Optimize approach

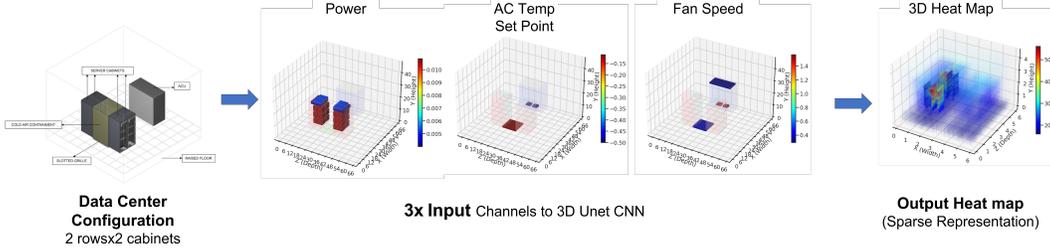


Figure 2: 3D CNN: Input and Output for Data Center Configuration

ground truth heatmaps, each of which we denote as a tensor $\mathbf{y} \in \mathbb{R}^{1 \times D \times H \times W}$. We denote an element within the output tensor \mathbf{y} as y_{ijk} and a set of output tensors as $\mathbf{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_N\}$. To create these ground truth heatmaps, denoted as tensors \mathbf{y} , we need dense targets for training. The original heatmaps obtained from 6SigmaDCX simulations are sparse and irregularly sampled in 3D space. To densify them, we use nearest interpolation on a voxel grid, as linear and cubic interpolations become computationally challenging with irregular data. To enhance the realism of the interpolated heatmaps, we apply a Gaussian filter with $\sigma = 0.5$ and a radius of 4σ . Finally, we constrain the temperature values to a range between 16°C and 50°C , addressing anomalies in less than 1% of the data.

Data Preprocessing and Training: Data for the 3D CNN model is normalized and shuffled each training epoch. For model parameter updates, we used the Adam optimizer with the smooth L1 loss to prevent large gradient updates across a diverse data set and make it less susceptible to outliers.

3D CNN Model Architectures: We evaluated our approach using four 3D CNN architectures. 1) **3D U-Net** Çiçek et al. (2016): An adaptation of U-Net Ronneberger et al. (2015) designed for 3D data, such as volumetric data. 2) **3D Residual U-Net** Lee et al. (2017): A modification of the 3D U-Net which incorporates residual convolutional blocks to support deeper networks. 3) **V-Net** Milletari et al. (2016): Originally introduced for volumetric medical image segmentation, it contains modifications like parametric ReLU activations and in-block residual connections. 4) **SegNet** Badrinarayanan et al. (2017): A 3D adaptation of the deep encoder-decoder structure of SegNet, which emphasizes memory efficiency with pooling indices.

3 Results

Model accuracy evaluation metrics : To evaluate the effectiveness of our approach, we consider several metrics: inference time, mean-squared-error (MSE), mean maximum absolute error (AE), mean AE, top- t AE, and 3D structural similarity (SSIM). Top- t AE is especially useful for evaluating how well our approach models hot spots in data centers. It is defined below:

$$\text{Top-}t \text{ AE}(\mathbf{Y}, \hat{\mathbf{Y}}) = \sum_{n=1}^N \sum_{i,j,k \in \text{top-}t(\mathbf{y}_n)} \frac{|y_{nijk} - \hat{y}_{nijk}|}{Nt} \quad (1)$$

The function $\text{top-}t(\cdot)$ returns the 3D indices of the largest t elements of its tensor argument, i.e., the hottest t temperature locations of a ground truth heatmap. In our evaluation, we set t to 10% of the number of heatmap elements. The inference time is the time to predict a single sample after the model and data are loaded into GPU memory.

Workload Settings for Data Generation: We collected data sets from CFD simulations across different parameter settings. In one case, we considered server utilization sampled uniformly from a predefined range: 25% to 95%. This simulated a normal use-case scenario. In certain downstream applications, the surrogate model may be evaluated using utilization that has wide variance, with one part of the data center having extremely low utilization while the other sections have extremely high utilization leading to extreme hotspots. In this regard, we considered a data generation with similar server utilization characteristics. Finally, we considered a more grid-based sampling to generate server utilization values. This was to create models that capture behaviors across a wide spectrum of values and sample data from all parts of the n -dimensional hypercube. We considered cases where the CFDs were generated from both no-containment and cold-containment scenarios. In total, there were four parameter settings for generating the data which are highlighted in Table 2.

Workload Settings	Model	Inf. Time (ms)	MSE	Top- t AE (C)	3D SSIM
Grid CPU Utilization	Res. U-Net	41.6	0.00118	1.68	0.9459
	U-Net	40.6	0.00018	0.53	0.9798
	V-Net	35.6	0.00312	2.35	0.8576
	SegNet	15.2	0.00559	5.26	0.8197
Grid CPU Utilization and Cold Aisle Containment	Res. U-Net	41.6	0.00011	0.43	0.9939
	U-Net	40.6	0.00009	0.44	0.9911
	V-Net	35.6	0.00176	1.84	0.8962
	SegNet	15.2	0.00315	4.73	0.9220

Table 1: Evaluation (Speed and Accuracy) metrics on test set. The inference time is computed on a V100 GPU.

Workload Settings	Model	Inf. Time (ms)	MSE	Mean AE (C)	Max AE (C)	Top- t AE (C)	3D SSIM
Uniform Utilization	Res. U-Net	41.6	0.00086	0.87	9.99	1.29	0.9522
	U-Net	40.6	0.00094	0.93	11.0	1.34	0.9437
	V-Net	35.6	0.00112	1.44	12.8	2.32	0.9112
	SegNet	15.2	0.00432	2.06	14.7	6.27	0.8434
Extreme Utilization	Res. U-Net	41.6	0.00010	0.34	4.65	0.37	0.9772
	U-Net	40.6	0.00031	0.57	9.90	0.62	0.9565
	V-Net	35.6	0.00048	0.72	6.66	0.79	0.9343
	SegNet	15.2	0.00137	1.25	8.83	2.74	0.9028
Grid CPU Utilization	Res. U-Net	41.6	0.00118	0.84	8.90	1.68	0.9459
	U-Net	40.6	0.00018	0.40	7.32	0.53	0.9798
	V-Net	35.6	0.00312	1.93	14.5	2.35	0.8576
	SegNet	15.2	0.00559	2.27	19.7	5.26	0.8197
Grid CPU Utilization and Cold Aisle Containment	Res. U-Net	41.6	0.00011	0.21	4.47	0.43	0.9939
	U-Net	40.6	0.00009	0.23	6.39	0.44	0.9911
	V-Net	35.6	0.00176	1.30	9.40	1.84	0.8962
	SegNet	15.2	0.00315	1.11	20.7	4.73	0.9220

Table 2: Evaluation (Speed and Accuracy) metrics on test set. The inference time is computed on a V100 GPU.

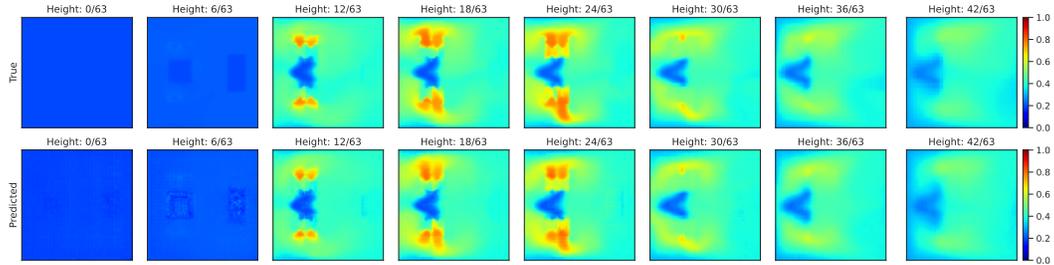


Figure 3: Matrix representation across different slices of the data center room. True outputs are in the top row, and model predictions the bottom row. Each column represents a slice at varying heights.

3D CNN Surrogate Results : When evaluated across various data center configurations and workloads, the Residual U-Net consistently delivered impressive results, notably achieving an MSE of 0.00002 and a mean AE of 0.14°C in certain scenarios. As a CFD surrogate, it achieved a significant speedup of over, 2800 times while maintaining reasonably effective accuracy to be used for downstream applications in workload distribution. Figure 3 further highlights an example of the faithful reproduction of the temperature distribution across 2D slices at different heights of the data center.

Genetic Algorithm Results : Using the genetic algorithm for workload optimization in data centers yielded considerable performance enhancements, notably reducing hotspots by approximately 7.70%. This algorithm converged efficiently, typically within 25 iterations (Figure 1 right). By integrating the 3D CNN surrogate model, we achieved over a 99% reduction in optimization time, showcasing the model’s efficiency for real-time data center optimizations.

Sustainability/Energy Results : Upon applying our optimized workload configurations derived from the genetic algorithm, we noticed a significant improvement in sustainability and energy efficiency. With an optimal workload distribution, the HVAC cooling energy consumption rates dropped by an average of 2.5%, translating to cost savings and a notable carbon footprint reduction. Moreover, the enhanced temperature regulation may lead to prolonged equipment lifespans and reduced e-waste and

manufacturing carbon footprint, emphasizing the holistic benefits of our approach. In summary, our results underscore the dual benefits of financial savings and environmental impact through optimized configurations.

4 Conclusion and Future Work

In our research, we developed a detailed methodology for modeling data center temperature distribution as a function of the geometric arrangement of the server workload. The resulting U-Net architectures showcased their efficacy as a surrogate model for 3D CNN in predicting and optimizing thermal configurations in data centers. Their impressive performance, marked by lower error metrics, becomes more pronounced when paired with a genetic optimization algorithm. This combination not only optimized server workload distribution but also significantly reduced optimization time, emphasizing its potential for real-time data center applications.

Our optimization efforts led to substantial improvements in energy efficiency and sustainability, evident from decreased energy usage and a reduced carbon footprint. Furthermore, prolonging equipment lifespans cuts down on e-waste, demonstrating this approach's combined economic and environmental advantages.

We plan to integrate the CFD surrogate into a digital twin for holistic multi-element and multi-objective real-time control to work towards the sustainability goal for lowering the carbon footprint.

References

- S. Ilager, K. Ramamohanarao, R. Buyya, Thermal prediction for efficient energy management of clouds using machine learning, *IEEE Transactions on Parallel and Distributed Systems* 32 (2020) 1044–1056.
- J. Lin, W. Lin, W. Lin, J. Wang, H. Jiang, Thermal prediction for air-cooled data center using data driven-based model, *Applied Thermal Engineering* 217 (2022) 119207. URL: <https://www.sciencedirect.com/science/article/pii/S1359431122011395>. doi:<https://doi.org/10.1016/j.applthermaleng.2022.119207>.
- S.-Q. Jin, N. Li, F. Bai, Y.-J. Chen, X.-Y. Feng, H.-W. Li, X.-M. Gong, W.-Q. Tao, Data-driven model reduction for fast temperature prediction in a multi-variable data center, *International Communications in Heat and Mass Transfer* 142 (2023) 106645. URL: <https://www.sciencedirect.com/science/article/pii/S0735193323000349>. doi:<https://doi.org/10.1016/j.icheatmasstransfer.2023.106645>.
- S. Fresca, A. Manzoni, Pod-dl-rom: Enhancing deep learning-based reduced order models for nonlinear parametrized pdes by proper orthogonal decomposition, *Computer Methods in Applied Mechanics and Engineering* 388 (2022) 114181.
- X. Chen, X. Zhao, Z. Gong, J. Zhang, W. Zhou, X. Chen, W. Yao, A deep neural network surrogate modeling benchmark for temperature field prediction of heat source layout, *Sci. China Phys. Mech. Astron.* 64 (2021) 1–30. doi:[10.1007/s11433-021-1755-6](https://doi.org/10.1007/s11433-021-1755-6).
- 6SigmaDCX | Future Facilities, CFD simulation tool, <https://www.futurefacilities.com/resources/videos/products/introducing-6sigmadcx/>, 2023. [Accessed 11-08-2023].
- Ö. Çiçek, A. Abdulkadir, S. S. Lienkamp, T. Brox, O. Ronneberger, 3d u-net: Learning dense volumetric segmentation from sparse annotation, in: *International Conference on Medical Image Computing and Computer-Assisted Intervention*, Springer, 2016, pp. 424–432. URL: <https://api.semanticscholar.org/CorpusID:2164893>.
- O. Ronneberger, P. Fischer, T. Brox, U-net: Convolutional networks for biomedical image segmentation, in: *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18*, Springer, 2015, pp. 234–241.

- K. Lee, J. Zung, P. Li, V. Jain, H. S. Seung, Superhuman accuracy on the snemi3d connectomics challenge, arXiv preprint arXiv:1706.00120 (2017).
- F. Milletari, N. Navab, S.-A. Ahmadi, V-net: Fully convolutional neural networks for volumetric medical image segmentation, in: 2016 Fourth International Conference on 3D Vision (3DV), IEEE, 2016. URL: <https://doi.org/10.1109/3dv.2016.79>. doi:10.1109/3dv.2016.79.
- V. Badrinarayanan, A. Kendall, R. Cipolla, SegNet: A deep convolutional encoder-decoder architecture for image segmentation, IEEE Transactions on Pattern Analysis and Machine Intelligence 39 (2017) 2481–2495. URL: <https://doi.org/10.1109/tpami.2016.2644615>. doi:10.1109/tpami.2016.2644615.