

Brain-Inspired Hyperdimensional Computing: How Thermal-Friendly for Edge Computing?

Paul R. Genssler, Austin Vas, and Hussam Amrouch

Abstract—Brain-inspired hyperdimensional computing (HDC) is an emerging machine learning (ML) methods. It is based on large vectors of binary or bipolar symbols and a few simple mathematical operations. The promise of HDC is a highly efficient implementation for embedded systems like wearables. While fast implementations have been presented, other constraints have not been considered for edge computing. In this work, we aim at answering how thermal-friendly HDC for edge computing is. Devices like smartwatches, smart glasses, or even mobile systems have a restrictive cooling budget due to their limited volume. Although HDC operations are simple, the vectors are large, resulting in a high number of CPU operations and thus a heavy load on the entire system potentially causing temperature violations. In this work, the impact of HDC on the chip's temperature is investigated for the first time. We measure the temperature and power consumption of a commercial embedded system and compare HDC with conventional CNN. We reveal that HDC causes up to 6.8 °C higher temperatures and leads to up to 47 % more CPU throttling. Even when both HDC and CNN aim for the same throughput (i.e., perform a similar number of classifications per second), HDC still causes higher on-chip temperatures due to the larger power consumption.

Index Terms—Hyperdimensional computing, Convolutional neural network, Embedded system, Edge computing, Temperature.

I. INTRODUCTION

WEARABLE devices are an emergent type of embedded system. Smartwatches, fitness trackers, e-textiles, or smart glasses are prominent examples. They detect seizures, measure the heart rate, categorize sport activities, augment reality to help the user, or enable wireless payments, among others. Developers of such devices are trending towards embedded intelligence for more independence from phones and also for immediate feedback. The required machine learning (ML) and deep learning methods typically involve heavy computations and require relatively large storage for the models. However, in embedded systems, resources such as processing power, memory, energy, and cooling are tightly constrained.

Brain-inspired hyperdimensional computing (HDC) has emerged as a promising alternative to the established ML methods. Similar or better inference accuracy compared to deep neural networks (NNs) has been demonstrated for speech recognition [6], EEG-based seizure detection [1], EMG-based gesture detection [7], or wafer map defect pattern classification [4], among others [3]. While layers of neurons are the base in NNs, large vectors with dimensions in the thousands, hypervectors, form the basic blocks in HDC. Their components can be

simple bits, bipolar values, or more complex numbers and they are all independent of each other. Hence, parallelization is easy and consequently, computationally efficient implementations have been proposed for ASICs [8], FPGAs [9], and embedded devices [1], [7].

Three basic operations like dot product, component-wise majority, and multiplication are used to encode the real-world data into non-binary hypervectors. Compared to the floating-point matrix multiplications required by NNs, such operations are lightweight, i.e., computationally less expensive. Therefore, HDC has been proposed as a promising ML method for embedded systems and smart things [1], [3], [6]–[9]. Other advantages include the robustness against noise in input data, e.g., from unreliable cheap sensors, or in the model's memory as well as the ability to learn from very little data [1]. To realize the fast product cycle demanded by consumers, off-the-shelf hardware and software implementations are preferred by developers for their short time to market. Highly specialized hardware like ASICs achieve better computational performance and higher energy efficiency. Nevertheless, HDC remains a demanding application due to the large size of the hypervectors, data movement, and complex encoding.

In addition to memory, time, and power constraints face embedded systems, wearables in particular, also a temperature constraint. Such systems cannot be actively cooled because fans draw additional power, contain error-prone moving parts, and require space that is not available, e.g., in smart glasses. Passive cooling is limited on one hand by the overall size of the system and thus its thermal capacity. On the other hand, if and how it has contact with human skin can drastically limit the maximum sustainable temperature. *The impact of HDC on the on-chip temperatures has not yet been evaluated.*

Our main contribution: We are the first to investigate if HDC conforms with the thermal constraint that an embedded system imposes. If it causes temperature violations due to the demanding operations, then the promise as an alternative ML method for off-the-shelf embedded systems needs to be revisited. We employ an optimized state-of-the-art implementation of HDC [5] on a commercial Raspberry Pi 4 to classify images from the Fashion-MNIST data set [10], a common task, e.g., in smart glasses. We measure the temperature of the CPU and the power consumption of the whole system under realistic workloads. HDC is compared against a convolutional neural network (CNN) during training and inference. We reveal that only HDC causes CPU throttling due to higher temperatures but still trains faster. In a scenario where both methods perform equally fast, HDC consumes more energy and causes higher on-chip temperatures.

Paul R. Genssler, Austin Vas, and Hussam Amrouch are with the Chair of Semiconductor Test and Reliability (STAR), University of Stuttgart, Stuttgart 70569, Germany. E-mail: {genssler, amrouch}@iti.uni-stuttgart.de.

arXiv:2204.03739v1 [cs.LG] 5 Apr 2022

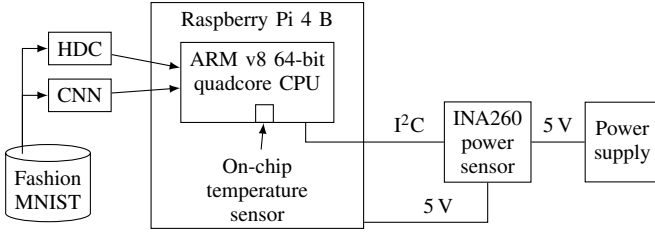


Fig. 1. Experimental setup with Fashion MNIST as a realistic workload.

II. PRELIMINARIES

The concepts of brain-inspired HDC can be implemented with different underlying vector symbolic architectures, e.g., with vectors of simple bits, bipolar values, integers, or real numbers. Depending on the type, different mathematical operations realize the basic operations of binding, bundling, permutation, and similarity score. In this work, real number hypervectors are employed because they offer the highest inference accuracy. A drawback is their higher computational cost compared to binary implementations.

To improve the inference accuracy, the concept of retraining has been proposed [6]. After an initial training cycle and class hypervector generation, the training data itself is used for inference. If a sample is incorrectly classified, then the class hypervector is adjusted to be more similar to the sample. The process is repeated. Each iteration is referred to as an epoch.

Temperature can be a limiting factor in embedded devices. It correlates with the amount of processing performed by the compute units, which can be CPUs, GPUs, or other specialized ones, like tensor processing units. In addition, memory accesses consume comparatively large amounts of electric energy, which translates again into thermal energy. Both sources of heat have to be controlled to achieve an overall lower system-wide temperature. Dynamic voltage and frequency scaling (DVFS) is one of the most important tools to manage the temperature. By reducing the clock frequency, the voltage can be reduced, having a quadratic impact on the dynamic power consumption.

III. EXPERIMENTAL SETUP

A Raspberry Pi 4 Model B with 8 GB RAM is employed in the experiments. The Raspberry Pi OS (32-bit Linux) uses kernel version 5.10.17-v7l and the firmware from May 27, 2021. The nominal CPU frequency is 1500 MHz and starts throttling once its temperature exceeds 80 °C consistently (“soft max”). In the experiments, this throttling prevents the temperature from reaching the maximum of 85 °C. Temperature is measured with the on-chip sensor of the Broadcom 2711 system on chip (SoC) and at idle is on average 45 °C. The ambient temperature is kept constant at about 20 °C in all experiments for fair comparisons. The board is not covered, in a case, nor is a heat sink attached. Power measurements are done with an intermediary INA260 sensor between the USB power adapter output and the board, as shown in Fig. 1. The measurements are collected by the Pi itself via I²C. Idle power consumption is about 3.5 W with the CPU frequency throttled down to 600 MHz so save power.

Image classification is a demanding task for embedded systems. Each sample contains a lot of data that has to

TABLE I
CNN ARCHITECTURES EMPLOYED IN THE EXPERIMENTS.

Layer type	Output shape	Parameters	Simplified CNN
Conv2D	(28, 28, 64)	320	(22, 22, 6)
MaxPooling2D	(14, 14, 64)	0	(11, 11, 6)
Dropout	(14, 14, 64)	0	(11, 11, 6)
Conv2D	(14, 14, 32)	8224	removed
MaxPooling2D	(7, 7, 32)	0	removed
Dropout	(7, 7, 32)	0	removed
Flatten	(1568)	0	(726)
Dense	(256)	401 664	(20)
Dropout	(256)	0	(20)
Dense	(10)	2570	(10)
Total trainable parameters		412 778	15050
Inference accuracy (70 epochs)		92.7 %	85.8 %

TABLE II
INFERENCE ACCURACY FOR DIFFERENT DIMENSIONS WITH HDC AND A NINE DEGREE POLYNOMIAL SVM FOR COMPARISON.

Dimension	500	2000	4000	10 000	SVM
ISOLET	92.2	93.2	93.4	93.6	96.3
MNIST	91.2	94.4	94.8	95.1	98.4
FMNIST	82.3	84.7	86.0	86.7	89.4

be processed. In this work, the Fashion-MNIST data set is employed as an example [10]. It comprises 60 000 training and 10 000 testing samples of ten different classes of clothing. All images are grayscale and have a resolution of 28x28 pixels.

OnlineHD [5] is a state-of-the-art HDC implementation using bipolar hypervectors. Version 0.1.2 is built on top of Python’s PyTorch framework and provides efficient implementations for the underlying mathematical operations during encoding. The similarity computation is written in C++.

Inspired by the visual mechanism in the brain, CNNs were proposed as a special form of a neural network with convolution layers, sub-sampling layers, and completely connected layers. Two CNN architectures are described in Tab. I and implemented with TensorFlow 1.14.0 and Python 3.7.3. The larger CNN is used in the experiments, while the simplified CNN is discussed in Section IV-C. Neither model utilizes the GPU during training or inference, all processing is done on the CPU.

IV. EXPERIMENTAL RESULTS

All tasks, such as training and inference, are continuously repeated for four hours. This allows the system to accumulate heat and to minimize environmental influences. The results presented in the figures are short representative segments from these long experiments. The reported numbers are averaged over the whole experiment.

A. Inference Accuracy

Each model, CNN and HDC, is trained on a server machine with no memory, temperature, or power constraint. Training the models longer does not increase the energy consumption or the thermal emissions of the deployed model. During training, only the neurons’ parameters or the hypervectors are changed, not the CNN architecture itself or the dimension. Therefore,

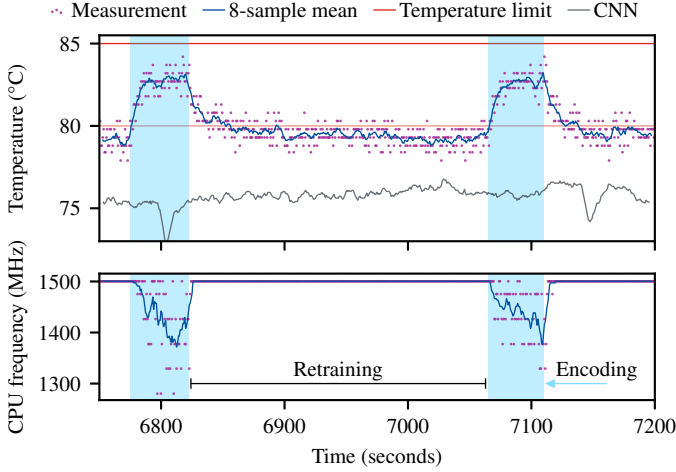


Fig. 2. Encoding during repeated training in HDC causes CPU throttling. Despite that, HDC is still faster in training the model.

fully trained models can be deployed to the target system. After training for 70 epochs, the CNN achieves an inference accuracy of 92.7% on the test set. For HDC, the various dimensions are explored and summarized in Tab. II, from 500 up to 10000 achieving an inference accuracy of 82.3% and 86.7%, respectively. Even though this accuracy increases with the dimension, anything above 4000 (86.0%) would have a diminishing return for an embedded system. It increases energy consumption and inference latency. Hence, a dimension of 4000 is selected for the remaining experiments.

B. Training Analysis on the Embedded System

The training of an ML model involves lots of data and computations. Hence, it is usually performed in the cloud and only the fully trained model is deployed to the target system. However, for personalized applications, like seizure detection or voice recognition, the model could be trained further to improve inference performance with the user. In this experiment, the model was trained from scratch to evaluate the overall impact of this type of workload on temperature. The data set is loaded once into memory and the training is repeated to simulate a larger data set. Re-/training epochs are reduced to one and 20, resulting in a reduced inference accuracy of 86.0% and 81.9% for the CNN and HDC, respectively. The CNN does not exceed 79°C with 75.6°C on average as shown in Fig. 2. In contrast, HDC reaches temperatures of up to 84.2°C, 6.7°C more than CNN on average. HDC is frequently close to the maximum temperature of 85°C during the encoding of the data set. Hence, DVFS is employed by the SoC, reducing the clock frequency up to 47% (800 MHz). During retraining, the temperature drops to about 79°C. No thermal buffers are created for the next encoding cycle. The power consumption is similar for both models with 5.2 W on average. Therefore, the operations used by HDC cause more heat compared to the CNN. Yet, despite the CPU throttling, HDC completes a training cycle 14.9% faster as listed in Tab. III. However, inference accuracy on the test set is 4% lower with HDC.

TABLE III
AVERAGE NUMBER OF IMAGES PROCESSED PER SECOND.

Scenario	Training	Inference	Inference
Temperature Limit	85 °C	85 °C	60 °C
HDC	206.3	1000.5	382.0
CNN	179.5	322.8	373.0

C. Continuous Inference Analysis

The primary task of an ML model in an embedded system is the inference of incoming samples. Smart glasses could analyze the live camera feed to overlay directions for an augmented reality experience. A smartphone listens for keywords to activate the voice assistant. To mimic a similar scenario, the inference of the test data set is continuously repeated.

HDC processes 1000.5 images per second, which is 3.1x more than CNN because of its less complex operations. It also utilizes the CPU to its full extent, whereas the CPU load during CNN inference is only about 75%. However, because of the more complex operations, the CNN consumes 5.87 W compared to HDC with 5.44 W, on average. Although power consumption is 8% higher with a CNN, the temperature is 2.1°C lower on average. Furthermore, the CNN exceeds the first temperature limit of 80°C for 7.2% of the total time compared to 47.5% with HDC. The trend in the data for CNN suggests an increase in temperature above 80°C if the data set was larger. However, the setup time between two inference runs reduces CPU load and temperature. During this setup period, book-keeping tasks are done and the inference accuracy calculated. The same applies to HDC, although the temperature increases faster and could reach 85°C, triggering DVFS. The measured CPU clock frequency reduction is not because of over-temperature but enabled by the reduced CPU load during setup.

A simplified CNN has been tested as well, its structure is described in Tab. I. Due to its reduced complexity, the number of processed images increase to 2427.5 outperforming HDC. Power consumption remains at the level of 5.3 W.

D. Inference Analysis under Temperature Constraint

A Raspberry Pi has an upper temperature limit of 85°C, an integrated heat spreader made of metal, and a large PCB to dissipate the thermal energy. Smaller edge devices, embedded systems, and wearables do not have such high temperature limits or thermal mass. To simulate these tighter temperature constraints with a Raspberry Pi, its clock frequency is throttled if a temperature threshold is crossed. In this experiment, the threshold is set to 60°C. If the devices touch the skin, such high temperatures could harm the user. A small script checks the temperature every second and sets the CPU frequency to 600 MHz in the case of overheating. As shown in Fig. 3, both models exceed the limit even with a throttled CPU clock. The number of processed images is reduced from 1001 (no limit, Section IV-C) to 382 for HDC and increases for CNN from 323 to 373. HDC has a 22% higher CPU load and consumes 7% more energy. This results in a 2°C higher temperature on average.

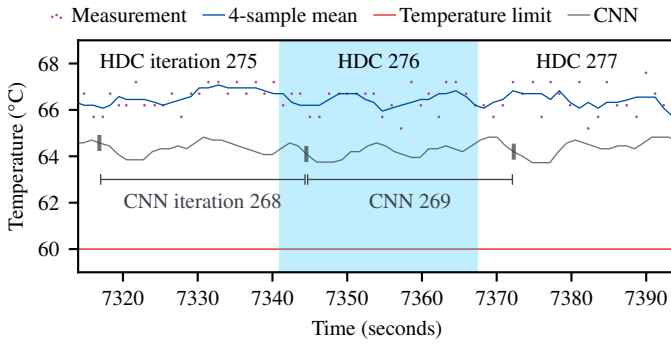


Fig. 3. Inference with a temperature constraints of 60°C. Both models cause permanent CPU throttling to 600 MHz. In addition, the temperature with HDC is 2°C higher but both classify a similar number of images per second.

E. Hand-crafted Encoding for ISOLET

To provide an additional perspective, the ISOLET dataset [2] of voice samples is analyzed. The first model is OnlineHD, which achieves an accuracy of 94.7% with a dimension of 2048. The second model is based on the C-Code implementation of HD-Lib¹, uses binary vectors of dimension 4096 and a record encoding to achieve an accuracy of 87.1% (without retraining). The third model is a KNN classifier with $n = 9$ and 92.1% accuracy implemented with the Scikit-learn library.

OnlineHD, HD-Lib, and KNN process 2369.5, 537.7, and 390.1 samples/s, respectively. OnlineHD also consumes the most energy (5.9 W) and generates the most heat (82.8°C) reducing the clock frequency to 1405 MHz on average. The HD-Lib-based hand-crafted encoding consumes slightly less power (5.7 W) and consequently generates less heat (82.1°C) but still throttles the CPU to 1431 MHz on average. The KNN does not utilize all CPU cores fully and the average temperature is about 74.9°C, not throttling the CPU and consuming less power (5.5 W).

V. DISCUSSION

① HDC performs training and inference faster than CNN. However, under temperature constraints and thus reduced CPU frequency because of throttling, the throughput (images per second) is reduced more with HDC compared to CNN. The number of images processed per second decreases by 33% and 62% during training and inference, respectively. ② In contrast, with CNN, the throughput decreases 16% and even increases 16% during training and inference, respectively. These results suggest HDC is *compute-bound* since a reduction in frequency by 60% reduces the inference time by 62%. The CNN, to the contrary, is *memory-bound*. The CPU is stalled waiting for data. A reduced frequency reduces this wait time and the CPU utilization increases. ③ Nevertheless, HDC consistently causes higher temperatures than CNN. During inference, HDC exceeded the targeted 80°C almost half of the time, 6.6x more often than CNN, yet HDC performs 3.1x more inferences per second. However, even with reduced CPU frequency, HDC exceeds the temperature limit by 2°C more than CNN, while both have comparable inference speed. Because of its high

number of computations compared to CNN, HDC performs worse regarding temperature and power consumption in such low-power scenarios expected for wearable devices. ④ The efficiency of the CNN is improved by using a 64-bit OS. The number of images processed per second jumps to 747 while the power consumption is reduced to 5.0 W at the same time. Alternatively, quantizing the CNN to 8-bit integers doubles the processing speed. A simplified CNN, described in Tab. I, achieves after 70 epochs training the same accuracy as HDC but processes 2427.5 images/s at 5.3 W, although at 77.7°C. Hand-crafted HDC encoding on a different dataset suggests a similar pattern. With off-the-shelf hardware, HDC is outperformed by other methods in non-functional metrics, like power consumption and temperature.

VI. CONCLUSION

This is the first work to investigate the impact of HDC on the on-chip temperature of a CPU. We conducted several experiments on a commercial embedded system. Power consumption and temperature are measured while executing HDC and CNN algorithms for both training and inference. We unveil that HDC faces a temperature and power challenge in low-power systems such as wearables and hence it is less thermal-friendly in edge computing.

ACKNOWLEDGMENT

This research was supported by Advantest as part of the Graduate School “Intelligent Methods for Test and Reliability” (GS-IMTR) at the University of Stuttgart.

REFERENCES

- [1] A. Burrello, S. Benatti, K. A. Schindler, L. Benini, and A. Rahimi, “An ensemble of hyperdimensional classifiers: Hardware-friendly short-latency seizure detection with automatic iEEG electrode selection,” pp. 1–1.
- [2] D. Dua and C. Graff, “Uci machine learning repository,” 2017. [Online]. Available: <https://archive.ics.uci.edu/ml/datasets/ISOLET>
- [3] L. Ge and K. K. Parhi, “Classification using hyperdimensional computing: A review,” vol. 20, no. 2, p. 30–47.
- [4] P. R. Genssler and H. Amrouch, “Brain-inspired computing for wafer map defect pattern classification,” in *2021 IEEE International Test Conference (ITC)*, pp. 123–132.
- [5] A. Hernandez-Cane, N. Matsumoto, E. Ping, and M. Imani, “OnlineHD: Robust, efficient, and single-pass online learning using hyperdimensional system,” in *2021 Design, Automation Test in Europe Conference Exhibition (DATE)*, pp. 56–61.
- [6] M. Imani, D. Kong, A. Rahimi, and T. Rosing, “VoiceHD: Hyperdimensional computing for efficient speech recognition,” in *2017 IEEE Int. Conf. on Rebooting Computing (ICRC)*, pp. 1–8.
- [7] A. Moin, A. Zhou, A. Rahimi, A. Menon, S. Benatti, G. Alexandrov, S. Tamakloe, J. Ting, N. Yamamoto, Y. Khan, F. Burghardt, L. Benini, A. C. Arias, and J. M. Rabaey, “A wearable biosensing system with in-sensor adaptive machine learning for hand gesture recognition,” vol. 4, no. 1, pp. 54–63.
- [8] F. Montagna, A. Rahimi, S. Benatti, D. Rossi, and L. Benini, “Pulp-hd: Accelerating brain-inspired high-dimensional computing on a parallel ultra-low power platform,” in *2018 55th ACM/ESDA/IEEE Design Automation Conference (DAC)*, 2018, p. 1–6.
- [9] M. Schmuck, L. Benini, and A. Rahimi, “Hardware optimizations of dense binary hyperdimensional computing: Rematerialization of hypervectors, binarized bundling, and combinational associative memory,” *J. Emerg. Technol. Comput. Syst.*, vol. 15, no. 4, Oct 2019.
- [10] H. Xiao, K. Rasul, and R. Vollgraf. (2017) Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms.

¹<https://github.com/skurella/hdlib>