

HDTVQ-VAE: Binary Codebook for Hyperdimensional Latent Representations

Austin J. Bryant[✉], *Graduate Student Member, IEEE*, and Sercan Aygun[✉], *Member, IEEE*

Abstract—Hyperdimensional computing (HDC) has emerged as a promising paradigm offering lightweight yet powerful computing capabilities with inherent learning characteristics. By leveraging binary hyperdimensional vectors, HDC facilitates efficient and robust data processing, surpassing traditional machine learning (ML) approaches in terms of both speed and resilience. This letter addresses key challenges in HDC systems, particularly the conversion of data into the hyperdimensional domain and the integration of HDC with conventional ML frameworks. We propose a novel solution, the hyperdimensional vector quantized variational auto encoder (HDTVQ-VAE), which seamlessly merges binary encodings with codebook representations in ML systems. Our approach significantly reduces memory overhead while enhancing training by replacing traditional codebooks with binary $(-1, +1)$ counterparts. Leveraging this architecture, we demonstrate improved encoding-decoding procedures, producing high-quality images within acceptable peak signal-to-noise ratio (PSNR) ranges. Our work advances HDC by considering efficient ML system deployment to embedded systems.

Index Terms—Hyperdimensional computing (HDC), vector quantized variational auto encoder (VQ-VAE).

I. INTRODUCTION

HYPERDIMENSIONAL computing (HDC) has recently received significant attention from researchers due to its ability to provide lightweight yet efficient computing coupled with certain learning characteristics [1]. HDC offers both speed and robustness compared to conventional machine learning (ML), owing to its unconventional data representation. Utilizing high-dimensional symbolic binary vectors, known as *hypervectors*, HDC performs various ML tasks [2]. These hypervectors possess the inherent property of holographic information storage [3]. By employing binary $(-1, +1)$ hyperdimensional vectors, HDC achieves holographic information storage [4], encompassing multiple incoming data within a single binary vector, with lower-resource consumption than conventional binary radix systems. The dimensionality (D) of hypervectors typically ranges from $1K$ to $10K$, with vectors necessitating near orthogonality for unbiased symbol representation. This unique orthogonality property facilitates *bundling*, *binding*, and *permuting* operations for encoding purposes, enabling a single vector to represent multiple inputs [5]. Bundling hypervectors yields a hypervector most similar to all

addends. The binding operation results in a hypervector most dissimilar to the set of hypervectors bound, akin to coordinate-wise multiplication. Permutation, a circular shift, generates uncorrelated hypervectors, preserving orthogonality [6].

While HDC is a promising path toward lightweight ML applications, much consideration must be given to how data is converted to its symbolic representation. This letter focuses on helping streamline the encoding process using unsupervised learning. In contrast to conventional ML methods, which often employ continuous, real-valued vectors, such as those in Word2Vec [7], the HDC approach utilizes a discrete, symbolic representation. This enables robust, memory-efficient holographic information storage using hyperdimensional binary vectors, reducing the overhead associated with floating-point representations. This letter uses conventional ML to help select symbols for representing abstract data. To the best of our knowledge, for the first time in the literature, vector quantized variational auto encoder (VQ-VAE) architecture is decorated with an HDC approach for an efficient latent space representation without the loss of generality. The VQ-VAE codebook's discrete aspect aligns with standard encoding techniques for HDC, where a discrete symbol list or *codebook* is used. In addition, the *encoder*, *decoder*, and *quantization* process of the VQ-VAE aligns with our goal of learning to map abstract data to a set of symbols in an unsupervised manner.

Once the codebook is obtained, the classification task becomes the second component of the overall learning architecture. Traditional systems typically rely on conventional floating-point memory representation involving complex operations, such as *multiply-accumulate-activate*. However, in HDC-based classifiers, a single-pass learning strategy is employed, diverging from traditional ML systems. This study integrates a binary codebook trained using HDTVQ-VAE with the learning strategy of the HDC paradigm. The latent space elements undergo further processing through binding and bundling to create a classifier system. Results indicate that this proposed framework offers an easily deployable architecture with reduced memory usage while maintaining accuracy levels comparable to conventional approaches. This letter presents the following contributions: 1) integration of HD mechanics into VQ-VAE architecture via codebook replacement; 2) enable the use of downstream HDC applications that take advantage of HD latent space; 3) reduction of latent space by $32\times$ from the transition of float32 codebook to the binary codebook; and 4) reduction of overall model size due to codebook replacement.

II. BACKGROUND

A. Hyperdimensional Computing

HDC is an emerging computing paradigm gaining attention, particularly in developing ultralightweight learning systems. In comparison to conventional ML approaches, HDC offers

Manuscript received 12 August 2024; accepted 12 August 2024. Date of current version 26 November 2024. This work was supported in part by the National Science Foundation under Grant 2019511, and in part by the start-up fund provided by the University of Louisiana at Lafayette for Dr. Sercan Aygun. This manuscript was recommended for publication by A. Shrivastava. (Austin J. Bryant and Sercan Aygun contributed equally to this work.) (Corresponding author: Austin J. Bryant.)

The authors are with the School of Computing and Informatics, University of Louisiana at Lafayette, Lafayette, LA 70504 USA (e-mail: austin.bryant1@louisiana.edu; sercan.aygun@louisiana.edu).

Digital Object Identifier 10.1109/LES.2024.3443881

1943-0671 © 2024 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission.

See <https://www.ieee.org/publications/rights/index.html> for more information.

direct data processing without the need for additional feature engineering, supports single-pass learning, eliminates the requirement for error optimizations and backpropagation, and achieves a lightweight model thanks to unary processing capabilities. Encoding data into hypervectors and applying several logic operations results in the learned class of the training sample, unlike the learning steps (*forward pass* - *backward pass*) in traditional neural network systems, which require several optimization steps.

HDC excels in representing symbols, such as letters, image pixel positions, and 1-D signal timestamps, by leveraging orthogonal binary vectors. Each vector comprises D dimensions and consists of randomly occurring $+1$ and -1 values (interpreted as logic-1 and logic-0 in a hardware environment).

For symbol representation, since it is not a numerical value, a vector evaluation follows the unbiased probability point between 0 and 1. Assuming the middle point of the overall probability is 0.5 for each symbol vector representation, no bias exists between generated vectors. The probability presents having a logic value of 1 over the total vector size, thus having equal probability for each symbol.

On the other hand, scalar values, including 1-D signal amplitude and image pixel intensity, are also effectively encoded within HDC systems. The minimum and maximum scalar values within the space are mapped to $-1-1-1-1 \dots -1$ (minimum hypervector) and $+1+1+1+1 \dots +1$ (maximum hypervector), respectively. Any value within a range is represented by randomly flipping bits from -1 to $+1$ (or vice versa) based on the magnitude of the value [8].

Once vector generation is complete, the data undergoes subsequent encoding steps, such as binding, shifting, permutation, and bundling, tailored to the specific application. Applications primarily use two types of encoding: 1) n -gram-based and 2) record-based encoding. In symbol-based applications like language processing, n -gram operation involves handling n symbols (letters) alongside their corresponding hypervectors. Any n group of symbols are combined using shifting and permutation of their corresponding hypervectors [9]. For instance, $n = 3$ -gram operations on text processing take three subsets of letter hypervectors in a sentence, shift and permute each 3-group vector, and multiply them through the binding operation. Multiplications of $+1$ and -1 binary values are implemented using XOR operations in hardware, and the final result is another hypervector. The resultant hypervectors (n -gram vectors) from each n -gram subset are then cumulatively aggregated (achieved through population count operations in hardware) to produce a single vector from multiple n -gram inputs. This accumulation is a bundling operation and facilitates the holistic learning property of HDC. The final step involves binarization of the accumulated values back to binary form. This process, customized for each dataset class, generates a unified 1-D vector representing the respective class. Each new sample contributes to this resultant class hypervector by undergoing the same encoding steps [10]. Inference follows the same encoding steps outlined above; classification entails a similarity assessment between the test sample vector and the class hypervectors. A higher-similarity score (measured through metrics, such as cosine similarity, Hamming distance, dot product, etc.) indicates a closer match between the class and the test hypervector, yielding classification results.

B. Vector Quantized Variational Auto Encoder

Variational auto encoders (VAE) are used in many applications, such as generative image diffusion models, to produce

a compressed latent space that is computationally less expensive to perform operations. The VQ-VAE architecture differs from conventional VAEs in learning discrete compressed representations instead of continuous ones. It replaces the reparametrization trick with a quantization process, mapping features via distance measurements to the nearest symbol from a finite symbol list or codebook [11]. This letter is motivated by the shared goal of encoding for HDC and the VQ-VAE encoding process, mapping data to a discrete symbolic representation. Fig. 1(a) illustrates conventional VQ-VAE architecture.

Conventional VQ-VAE uses three main loss metrics to guide its training: 1) mean square error (MSE); 2) *Codebook Loss*; and 3) *Commitment Loss*. The MSE is calculated between the reconstructed and original data. Codebook loss and commitment loss are nearly identical. The only difference between them is which set of vectors is detached to force the other to be considered during backpropagation. In what follows, the losses in the conventional VQ-VAE are given:

$$\text{MSE Loss} = \frac{1}{N} \sum_{i=1}^N (x_{\text{recon},i} - x_i)^2$$

$$\text{Codebook Loss} = \text{MSE}(\text{quantized}, x_{\text{detach}})$$

$$\text{Commitment Loss} = \text{MSE}(\text{quantized}_{\text{detach}}, x)$$

III. HDVQ-VAE FRAMEWORK

The proposed architecture builds upon a well-known conventional ML module: the VQ-VAE, with an update of binary data representation to incorporate orthogonal vectors, resulting in the HDVQ-VAE. This adaptation offers a lightweight, easily deployable, and more efficient classifier utilizing the encoding dynamics of the HDC paradigm. Our architecture comprises *down blocks*, *up blocks*, and *quantization* modules. The down blocks consist of dual sequential modules, each comprising two convolution layers through which the input passes before undergoing addition and being fed through the *tanh* activation function. The first of these dual sequential modules features a preactivation *tanh* function situated between the convolution layers. Fig. 1(a) illustrates the conventional and Fig. 1(b) shows the proposed HDVQ-VAE architecture.

By replacing the convolution layers in the down block with convolutional transpose layers, we created our up block. The quantization module holds the codebook and is responsible for mapping latent vectors to this codebook. In our proposal, we replace the conventional codebook with a binary ($-1, +1$) representation for our architecture. This allows us to utilize hypervectors or sections of hypervectors as our codebook. We initialize these codebook vectors using pseudo-random sequences. In HDC, where orthogonality is essential for symbol representation, random sources provide this orthogonality.

Once the codebook is initialized, it remains fixed throughout the overall architecture without further updates, unlike a conventional VQ-VAE. This enables us to eliminate the codebook and commitment loss. We maintain the straight-through estimator to propagate gradients from one side of the quantization module to the other [11]. The codebook can comprise either hypervectors or sections of hypervectors. If the codebook consists of sections of hypervectors, then the latent space must contain multiple vectors that, when flattened, form a hypervector.

We have significantly streamlined the training process by eliminating the codebook and commitment loss in our

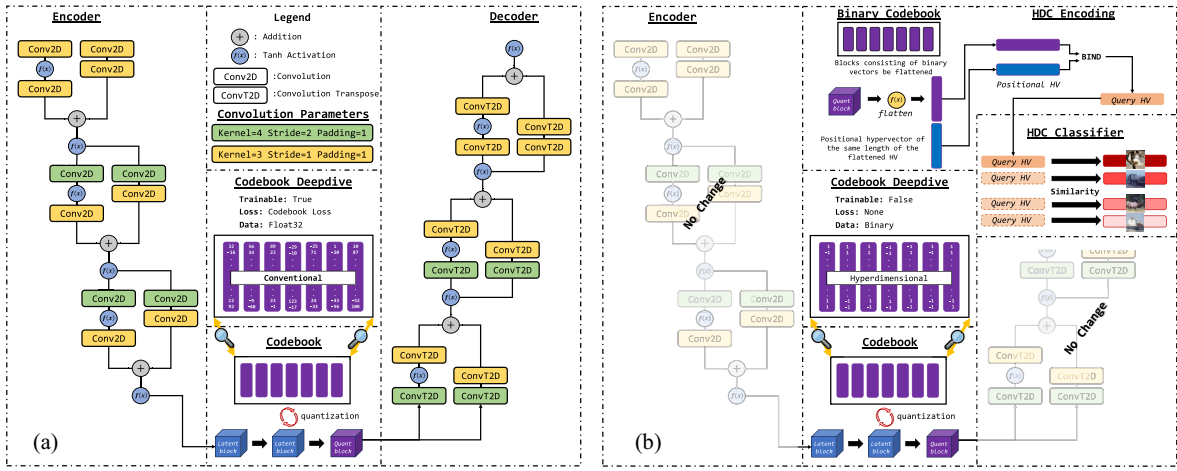


Fig. 1. Vector quantized variational auto encoder: (a) Conventional architecture fine-tuned for this letter and (b) HDVQ-VAE: Proposed architecture using HDC for latent spaces; the new classifier part is in HDC learning approach.

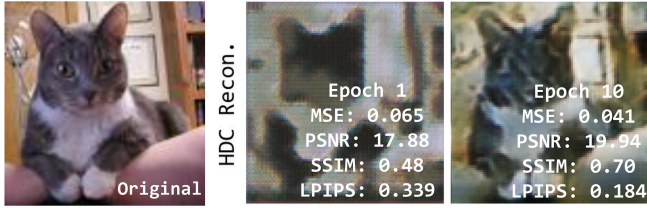


Fig. 2. Generative performance of the latent spaces in HDC domain.

modified VQ-VAE model. The static codebook does not require updates during training, making the codebook loss unnecessary. Similarly, the binary nature of the codebook diminished the benefits of a commitment loss. In addition to the conventional loss function, we have opted to create loss functions (L) from the well-known metrics to help guide the network: peak signal-to-noise ratio (PSNR) [12], structural similarity (SSIM), learned perceptual image patch similarity (LPIPS) [12], and perplexity (PPLX) [13]:

$$\text{PSNR}_L = 1.0 - \frac{\text{PSNR}(x_i, x_{\text{recon}})}{48} \quad \|\text{SSIM}_L = 1.0 - \text{SSIM}(x_i, x_{\text{recon}})$$

$$\text{LPIPS}_L = \frac{1}{N} \sum_{i=1}^N \text{LPIPS}(x_i, x_{\text{recon},i}) \quad \|\text{PPLX}_L = \frac{1}{\exp(-\sum p_i \log p_i)}$$

In our experiment, the training procedure uses the loss of the reconstructed images and their originals to guide the model toward producing meaningful latent representations. The training procedure is as in Algorithm 1.

In addition to utilizing HDC for binary latent space representation, the classifier component is also retained within the HDC domain to further leverage the latent space. Fig. 1(b) further illustrates the HDC classifier specifically designed for the codebook in the HDC domain. The process involves multiplying (binding) and accumulating flattened vectors with positional hypervectors for each corresponding class. The positional vectors coincide with our latent block's spatial dimensionality of 24×24 . Initializing our positional vectors involves assigning each spatial position a distinct binary. The distinctness of these vectors is ensured by the random assignment of bits to 0 or 1 with a 0.5 probability of either. Similarly to the quantized latent block, we flatten the positional vectors into a single hypervector that is then bound with the flattened quantized latent block. When evaluating (inference phase) incoming query samples, they are

Algorithm 1 HDVQ-VAE Training Procedure

- 1: **Initialize:**
 - 2: Model parameters use default Pytorch initialization
 - 3: ★codebook is initialized as a stack of nontrainable binary vectors
 - 4: **for each epoch do**
 - 5: **for each batch in training data do**
 - 6: **Encode:**
 - 7: Input data is $B \times 3 \times 96 \times 96$
 - 8: Result is latent block of size $B \times 128 \times 24 \times 24$
 - 9: **Quantize:**
 - 10: Compute distances to ★codebook vectors
 - 11: Assign each latent vector to the nearest ★codebook vector
 - 12: **Decode:** Generate reconstruction from quantized latent block
 - 13: **Calculate Losses * :**
 - 14: Compute Reconstruction, Perplexity, PSNR, SSIM, and LPIPS losses, then add all losses to calculate the total loss for the batch
 - 15: **Backpropagate:** Update model parameters based on total loss
 - 16: **end for**
 - 17: **end for**
- ★ The codebook consists of binary vectors (size of 128), unlike conventional VQ-VAEs. * Codebook loss is not used in HDVQ-VAE.

similarly encoded in the hypervector domain and compared to each previously generated class hypervector. The highest score determines the classification result. Consequently, this letter introduces two primary proposals to be tested in the subsequent section: 1) binary latent spaces and 2) codebook classification using HDC classifier.

IV. EXPERIMENTS AND RESULTS

In this section, we run two groups of experiments considering our two main contributions: 1) the performance of HD latent space for reconstruction and 2) latent space performance in the classification problem. Using the architecture in Fig. 1(b), we first consider the performance of the binary HD latent space. The images are constructed based on the STL10 dataset [14], and the construction performance of the HDC-based approach is in the acceptable performance range.

Our experiments use the STL10 dataset [14]. We chose the STL10 dataset for its larger $96 \times 96 \times 3$ ($W \times H \times C$) size compared to MNIST and CIFAR datasets. The larger size allowed us to test the reconstruction of higher-quality images where the codebook vectors may need to represent more details with a challenging problem.

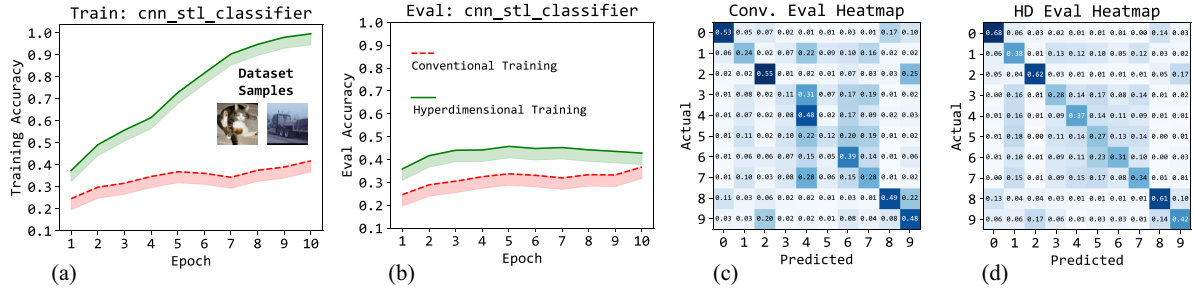


Fig. 3. CNN-based conventional classifier performance over the scalar latent space and HD latent space classification using the CNN model.

The additional losses focused on detail help guide the model in the selection of codebook vector to best pass detail information from the encoder to the decoder. The HDVQ-VAE can reconstruct images acceptably, as can be seen in Fig. 2.

Once trained, we can produce binary latent vectors using the encoder and binary codebook in Fig. 1(b). We use the latent vectors to test the classification abilities of the conventional method with respect to the convolutional neural network (CNN). The CNN classifier consists of three 2D convolutional layers, a flattened layer, and two fully connected layers. The convolutional layers have input-output channels of 128, 64, 32, and 8, each followed by a tanh activation. After flattening, the output passes through two fully connected layers, reducing from 72 to 32 features with ReLU activation, then to 10 features. We trained the model using a “reduce on plateau” scheduler with an initial learning rate of 0.003, a 0.5 scale factor, and a 0.01 threshold for cross-entropy loss. As seen in Fig. 3(a) and (b), the CNN classifier nearly doubled in evaluation accuracy while more than tripling in training accuracy. Fig. 3(c) and (d) report the testing phase confusion matrix for conventional and HD approaches, respectively.

For our HDC classification test, we reshape the binary latent vectors from our HDVQ-VAE, flattening them into a much higher-dimensional vector, which forms our hypervector. The flattening operation brings our latent representation from $B \times 128 \times 24 \times 24$ to $B \times 73\,728$ where B is our batch dimension, and 128 is the binary vector size. Our new hypervector is bound with a static positional hypervector initialized at the beginning of training. Once bound, we do a single pass to accumulate vectors for their respective classes before signing. As we begin training, we take similarity measurements between the bound hypervector ($latent \times position$) and our class vectors. We adjust the class vectors based on incorrect predictions. At the end of each epoch, we once again sign the class vectors.

While training an HDC classifier does not use backpropagation, it outperformed conventional classifier models on loss and accuracy scores. VQ version of VAE is more compliant with HDC than CNN due to the intricate learning of the quantization process. Not only is HDC computationally less expensive, but it is also lightweight, making it a suitable candidate for embedded systems.

The binarization of the codebook for our HDVQ-VAE brings memory savings and an overall reduction in the model size. Going from float32 to a binary representation shrinks the codebook size by $32\times$, and the information passed to the decoder is $32\times$ smaller than that of a conventional VQ-VAE. Despite this $32\times$ reduction in the information flow, comparing downstream latent classification, both the conventional and HDC classifiers perform better when using the latent of the HDVQ-VAE over the conventional VQ-VAE.

V. CONCLUSION

This study introduces a novel approach to enhance vector quantized variational autoencoders by incorporating hyperdimensional vector representations, aiming for a more efficient and lightweight solution in generative artificial intelligence (AI) and latent space classifications. By leveraging a binary hypervector space for latent representations, the codebook within the vector space can be utilized for image reconstructions or classifications. This research marks the inception of hyperdimensional computing (HDC) in machine learning (ML), suggesting its potential as a valuable addition to the conventional ML framework in future endeavors.

ACKNOWLEDGMENT

The open-source code for this letter can be found at <https://github.com/tMLEClab/HDVQ-VAE>.

REFERENCES

- [1] L. Ge and K. K. Parhi, “Classification using hyperdimensional computing: A review,” *IEEE Circuits Syst. Mag.*, vol. 20, no. 2, pp. 30–47, 2nd Quart., 2020.
- [2] S. Aygun, M. S. Moghadam, M. H. Najafi, and M. Imani, “Learning from hypervectors: A survey on hypervector encoding,” 2023, *arXiv:2308.00685*.
- [3] C. Yeung, Z. Zou, and M. Imani, “Generalized holographic reduced representations,” 2024, *arXiv:2405.09689*.
- [4] P. Kanerva, “Hyperdimensional computing: An introduction to computing in distributed representation with high-dimensional random vectors,” *Cogn. Comput.*, vol. 1, no. 2, pp. 139–159, Jun. 2009.
- [5] M. Stock et al., “Hyperdimensional computing: A fast, robust and interpretable paradigm for biological data,” 2024, *arXiv:2402.17572*.
- [6] A. Rahimi, P. Kanerva, and J. M. Rabaey, “A robust and energy-efficient classifier using brain-inspired hyperdimensional computing,” in *Proc. ISLPED*, 2016, 64–69.
- [7] A. G. Ayar, S. Aygun, M. H. Najafi, and M. Margala, “Word2HyperVec: From word embeddings to hypervectors for Hyperdimensional computing,” in *Proc. Great Lakes Symp. VLSI*, 2024, pp. 355–356.
- [8] T. Basaklar, Y. Tuncel, S. Y. Narayana, S. Gumussoy, and U. Y. Ogras, “Hypervector design for efficient hyperdimensional computing on edge devices,” 2021, *arXiv:2103.06709*.
- [9] A. Joshi, J. T. Halseth, and P. Kanerva, “Language geometry using random indexing,” in *Proc. 10th Int. Conf. Quantum Interact.*, 2017, pp. 265–274.
- [10] A. Hernández-Cano, N. Matsumoto, E. Ping, and M. Imani, “OnlineHD: Robust, efficient, and single-pass online learning using hyperdimensional system,” in *Proc. DATE*, 2021, pp. 56–61.
- [11] A. van den Oord, O. Vinyals, and K. Kavukcuoglu, “Neural discrete representation learning,” 2018, *arXiv:1711.00937*.
- [12] O. Keleş, M. A. Yilmaz, A. M. Tekalp, C. Korkmaz, and Z. Dogan, “On the computation of PSNR for a set of images or video,” 2021, *arXiv:2104.14868*.
- [13] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, “The unreasonable effectiveness of deep features as a perceptual metric,” 2018, *arXiv:1801.03924*.
- [14] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever, “Language models are unsupervised multitask learners,” *OpenAI Blog*, vol. 1, no. 8, p. 9, 2019.
- [15] A. Coates, A. Ng, and H. Lee, “An analysis of single-layer networks in unsupervised feature learning,” in *Proc. 14th Int. Conf. Artif. Intell. Stat.*, 2011, pp. 215–223.