

# Early Termination for Hyperdimensional Computing Using Inferential Statistics

Pu (Luke) Yi

lukeyi@stanford.edu  
Stanford University  
USA

Chae Young Lee

chae@stanford.edu  
Stanford University  
USA

Yifan Yang

yyang29@stanford.edu  
Stanford University  
USA

Sara Achour

sachour@stanford.edu  
Stanford University  
USA

## Abstract

Hyperdimensional Computing (HDC) is a brain-inspired, lightweight computing paradigm that has shown great potential for inference on the edge and on emerging hardware technologies, achieving state-of-the-art accuracy on certain classification tasks. HDC classifiers are inherently error resilient and support early termination of inference to approximate classification results. Practitioners have developed heuristic methods to terminate inference early for individual inputs, reducing the computation of inference at the cost of accuracy. These techniques lack statistical guarantees and may unacceptably degrade classification accuracy or terminate inference later than is needed to obtain an accuracy result.

We present OMEN, the first dynamic HDC optimizer that uses inferential statistics to terminate inference early while providing accuracy guarantees. To realize OMEN, we develop a statistical view of HDC that reframes HD computations as statistical sampling and testing tasks, enabling the use of statistical tests. We evaluate OMEN on 19 benchmark instantiations of four classification tasks. OMEN is computationally efficient, delivering up to 7.21–12.18× inference speed-ups over an unoptimized baseline while only incurring a 0.0–0.7% drop in accuracy. OMEN outperforms heuristic methods, achieving an additional 0.04–5.85× inference speed-up over the unoptimized baseline compared to heuristic methods while maintaining higher or comparable accuracy.

**CCS Concepts:** • **Hardware** → Emerging languages and compilers; Memory and dense storage; • **Software and its engineering** → Software notations and tools.

**Keywords:** program optimization, unconventional computing, hyperdimensional computing, emerging hardware technologies, edge machine learning

## ACM Reference Format:

Pu (Luke) Yi, Yifan Yang, Chae Young Lee, and Sara Achour. 2025. Early Termination for Hyperdimensional Computing Using Inferential Statistics. In *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 1 (ASPLOS '25)*, March 30–April 3, 2025, Rotterdam, Netherlands. ACM, New York, NY, USA, 17 pages. <https://doi.org/10.1145/XXXXXX.XXXXXX>

## 1 Introduction

Lightweight machine learning (ML) models have been effective in performing classification on the edge for biomedical, industrial, and environmental monitoring domains—examples include anomaly detection, wake word recognition, action prediction, and classification of signal patterns [22, 36]. In recent years, researchers and practitioners have developed ML classifiers that leverage hyperdimensional computing (HDC), an emerging computing paradigm, to perform inference on highly resource-constrained edge devices and error-prone emerging hardware platforms [6, 32, 38, 48, 65]. HDC classifiers are attractive targets for these platforms because their operators are resource efficient and amenable to optimization, the compute model is highly robust to errors, and inference results can be easily approximated.

### 1.1 HDC Classifiers and Early Termination

Hyperdimensional computing (HDC) is an emerging computing paradigm that uses high-dimensional vectors called hypervectors as the basic data type. In HDC, information is encoded in the distances between hypervectors. This representation is inherently error resilient and amenable to approximation. Data are encoded with hypervectors through the use of HD operators which compute hypervectors that

---

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).  
*ASPLOS '25, March 30–April 3, 2025, Rotterdam, Netherlands.*

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 979-8-4007-0698-1/25/03

<https://doi.org/10.1145/XXXXXX.XXXXXX>

are similar/dissimilar to their hypervector operands. Because all HD operators are element-wise operators and distances are approximable, HD computations are implementable as streaming computations that may be terminated early to obtain approximate results.

**Early Termination.** Due to the *fully distributed* property of HDC, the inference task is error-resilient and can be terminated at any time to obtain an approximation of the final classification result. HDC classifiers, therefore, support performance optimizations that terminate inference tasks early to reduce latency. In contrast, classical ML models typically deploy all-or-nothing inference algorithms in which the entire model must be applied to the inference input to produce a classification or offer coarse-grained early exit without guarantees on accuracy loss [23].

**Prior work.** Researchers have developed dynamic, heuristic approaches to terminate inference early in the runtime and methods to statically reduce the size of the hypervector [9, 30, 33, 52, 53, 61]. These methods lack statistical guarantees, and as a result may produce unacceptable accuracy degradations or process more hypervector elements than necessary. Researchers have also developed methods that minimize hypervector size while providing theoretical guarantees; however, these methods make independence assumptions about encoded information that do not hold for classification tasks [72].

## 1.2 Early Termination with OMEN

We present OMEN, a dynamic HDC optimizer that uses inferential statistics to terminate classification early with accuracy guarantees. OMEN deploys a termination algorithm inspired by Wald’s sequential probability ratio test (SPRT) [68], where new data are processed sequentially until a statistically significant conclusion is reached. OMEN offers the following benefits:

- **Fine-Grained Per-Input Early Termination.** OMEN identifies effective termination points for individual inputs at the dimension granularity. OMEN can therefore allocate compute resources based on the difficulties of inputs.
- **Accuracy Guarantees.** OMEN works with a user-provided confidence level that serves as an upper bound on the accuracy drop resulting from early termination. These accuracy guarantees hold even with hardware-induced errors.
- **Practically Deployable.** OMEN’s termination algorithm is optimized to introduce minimal memory and compute overheads and can be readily deployed on embedded systems. Furthermore, because OMEN requires little to no modification of the HDC training and inference algorithms, it can be readily integrated with other methods.

**Results.** OMEN is computationally efficient, delivering up to 7.21–12.18× inference speed-ups over an unoptimized baseline while only incurring a 0.0–0.7% drop in accuracy. Moreover, compared to heuristic early termination methods, OMEN achieves an additional 0.04–5.85× inference speed-up over the unoptimized baseline with higher or comparable accuracy

in 19 benchmarks. We also demonstrate that OMEN can be readily used without adaption to perform principled early termination in the presence of bit corruptions.

**Key Insight.** To realize OMEN, we develop a new statistical view of HDC that reframes HD computations as statistical sampling and testing problems. We relate the statistical exchangeable property to the *fully-distributed* property of hypervectors, which enables such a view, and show that this property holds for hypervectors constructed using simple HD-operator-based or advanced optimization-based algorithms.

## 1.3 Contributions

- **Statistical view of HDC** (Section 4): We present a statistical view of HDC. To the best of our knowledge, this is the first work that frames HDC classification as a statistical sampling and testing problem, allowing the application of statistical methods to HDC classifiers.
- **OMEN optimizer** (Section 5): We propose OMEN, the first dynamic early termination algorithm for HDC classifiers that uses inferential statistics to improve inference performance while providing accuracy guarantees.
- **Exchangeability** (Section 6): We relate the *fully distributed* property of HDC to statistical exchangeability and prove that hypervectors constructed using advanced training algorithms are exchangeable.
- **Evaluation** (Section 7): We empirically demonstrate that OMEN achieves significant performance improvements with minimal accuracy loss, and OMEN’s accuracy bounds hold robustly even in the presence of hardware errors.

## 2 Hyperdimensional Computing (HDC)

Hyperdimensional Computing (HDC), or Vector Symbolic Architectures (VSA), is a highly error-resilient, brain-inspired computing paradigm. The basic HDC datatype is a *hypervector*,  $V$ , a high-dimensional vector, and the value in each index  $V[i]$  is referred to as a *hypervector element*. The datatype of the hypervector elements depends on the variant of HDC used; we will target two variants in this work, BSC (Binary Spatter Code) and MAP (Multiply-Add-Permute). BSC uses dense binary hypervectors with elements  $V[i] \in \{0, 1\}$ , and MAP uses real-valued hypervectors with elements  $V[i] \in [-1, 1]$ . Table 1 summarizes the HD operator implementations.

**Distances.** In HDC, information is represented as relative distances between hypervectors. BSC uses the Hamming distance  $HD(V, V')$  between hypervectors and MAP uses the cosine  $\cos(V, V')$  or dot product  $\text{dot}(V, V')$  similarities between hypervectors. The size ( $N$ , number of elements) of the hypervector affects the resolution of encoded information and the granularity at which distances can be resolved.

**HD Codebooks.** To work with HDC, input datatypes, such as text, numbers, or images, must be encoded as a hypervector

Type	Hypervector	Random Generation	Distance Metric [ $dis(V, V')$ ]	Binding [ $V \otimes V'$ ]	Bundling [ $\oplus$ ]	Permutation [ $p^k(V)$ ]
BSC [37]	Dense Binary $V[i] \in \{0,1\}$	Bernoulli $V[i] \sim \text{Bern}(p=0.5)$	Hamming Distance $HD(V, V')$ $HD(V, V') = \frac{1}{N} \sum (V[i] \text{ xor } V'[i])$	XOR $V[i] \text{ xor } V'[i]$	Majority Vote $[(\sum_{j=1}^m V_j[i]) > \frac{m}{2}]$	Rotational Shift $V[(i+k) \bmod N]$
MAP [56]	Reals $V[i] \in [-1,1]$	Uniform $V[i] \sim U(-1,1)$	Cosine/Dot Product Similarity $dot(V, V') = \sum V[i] V'[i]$	Multiplication $V[i] V'[i]$	Normalized Addition $clamp(\sum V_j[i])$	Rotational Shift $V[(i+k) \bmod N]$

**Table 1.** BSC and MAP HDC Overview. Distance metrics compute a real value from two hypervectors. Other elementwise operators map hypervectors to hypervectors. *clamp* operation clamps the value to  $[-1,1]$  range.

representation. The basic building block of an HD information encoding are basis or code hypervectors which represent atomic symbols in the problem domain. These hypervectors are typically randomly generated, where each hypervector element  $V[i]$  is independently and identically distributed (*iid*).

**HD Operators.** HDC’s compute operators generate hypervectors with certain distance relationships to other hypervectors. The binding ( $V \otimes V'$ ) and permutation ( $p^k(V)$ ) operations produce hypervectors that are dissimilar (far) from the input hypervectors  $V$  and  $V'$ . The bundling operation ( $V_1 \oplus \dots V_m$ ) produces a hypervector that is similar (close) to each of the input hypervectors  $V_1, \dots, V_m$ . The implementation of HD operators depends on the HDC variant, as summarized in Table 1.

### 2.1 Example: Text Encoding

Consider the case where we want to encode the word **catch** as a hypervector. We first generate a random hypervector for each letter  $a \rightarrow V_a, \dots, z \rightarrow V_z$ . These letter hypervectors are far apart from one another. Using HD operators, we can then build hypervector representations of strings from the letter hypervectors. We note that the datatype of the generated hypervector  $V_{\text{catch}}$  depends on the HDC variant used. The HD operators determine the distance relationships between words; all the following three encodings encode the word **catch** but have very different distance properties:

**Binding-Based.** We encode the position of each letter by permuting the letter hypervector and then bind them together. The resulting hypervector embedding for **catch** will be dissimilar to words that share letters in the same position, such as **patch** and **match**:

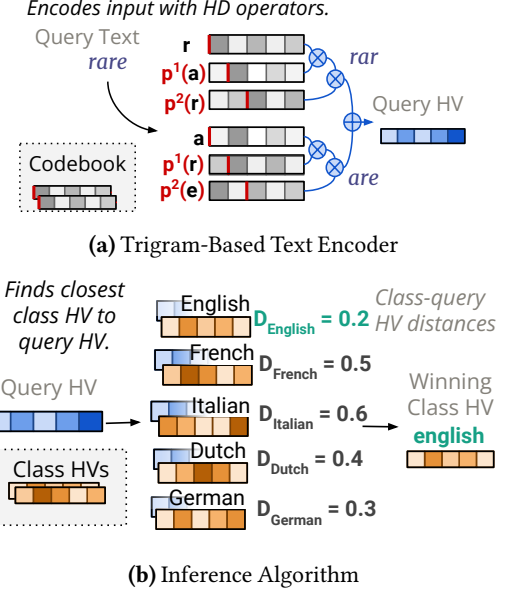
$$V_{\text{catch}}^{\text{bind}} = V_c \otimes p^1(V_a) \otimes p^2(V_t) \otimes p^3(V_c) \otimes p^4(V_h)$$

**Bundling-Based.** We encode the position of each letter by permuting the letter hypervector and then bundle them together. The resulting hypervector embedding for **catch** will be similar to words that share letters, such as **patch** and **match**:

$$V_{\text{catch}}^{\text{bund}} = V_c \oplus p^1(V_a) \oplus p^2(V_t) \oplus p^3(V_c) \oplus p^4(V_h)$$

**Trigram-Based.** We can build trigrams **cat**, **atc**, and **tch** with binding-based encoding and then bundle the trigrams together to build the word. With this encoding, words that share trigrams will be similar to each other.

$$V_{\text{catch}}^{\text{trig}} = V_{\text{cat}}^{\text{bind}} \oplus V_{\text{atc}}^{\text{bind}} \oplus V_{\text{tch}}^{\text{bind}}$$



**Figure 1.** HDC language classifier.

## 3 HDC Classification with OMEN

HDC has been used extensively for performing lightweight and error-resilient classification tasks, and has been demonstrated to achieve higher classification accuracy than other lightweight ML approaches on a range of edge inference tasks [26, 54]. OMEN enables early termination of inference tasks when a statistically significant result is reached, largely reducing the computational costs of inference. We next introduce OMEN using a language classification problem as a running example. Section 3.1 overviews the anatomy of HDC classifiers using the HDC language classifier. Section 3.2 introduces the challenges associated with reducing hypervector size during inference, and Section 3.3 introduces OMEN.

### 3.1 Anatomy of an HDC Classifier

Figure 1a and 1b presents the flow for an HDC language classifier. Given an input sentence, it encodes the sentence as a hypervector and classifies the sentence as one of the languages by computing the hypervector similarities. This classifier was previously applied to an in-memory computing substrate to enable language classification with high accuracy and very low energy consumption [38].

**Architecture.** The classifier is composed of an encoder (Figure 1a) that translates sentences to hypervectors. Hypervectors of training data are used to construct *class hypervectors*  $V_{\text{English}}$ ,  $V_{\text{French}}$ ,  $V_{\text{Dutch}}$ ,  $V_{\text{German}}$ , and  $V_{\text{Italian}}$  that capture the salient features of each language. In the above language classifier, the encoder uses a pre-generated codebook to translate each symbol (letters and punctuations) into a hypervector, and the trigram encoding method described in Section 2.1 to encode sentences by regarding the sentence as a large “word”.

**Inference.** During inference (Figure 1b), the classifier accepts a query sentence, translates it into a hypervector  $V_{\text{query}}$  using the encoder (Figure 1a), and then finds the class hypervector with the smallest distance to the query hypervector  $\text{argmin}_{\text{lang}}(\text{dis}(V_{\text{query}}, V_{\text{lang}}))$ .

**Training.** The HDC training algorithm computes class hypervectors from the training data, where each training input  $t_{i,l}$  is the  $i$ -th input with label  $l$ . Given a training dataset, the simplest HDC training algorithm first encodes each training input  $t_{i,l}$  as a hypervector  $V_{i,l}$  using the classifier’s encoder. The hypervectors are then grouped by label and then bundled together  $V_l = V_{1,l} \oplus V_{2,l} \dots$  to produce a class hypervector  $V_l$  for each label. More advanced training algorithms achieve higher accuracy by iteratively re-bundling misclassified inputs or learning hypervector representations. These methods are also supported by OMEN and discussed in Section 6.2.4. In this example, we use iterative rebundle training.

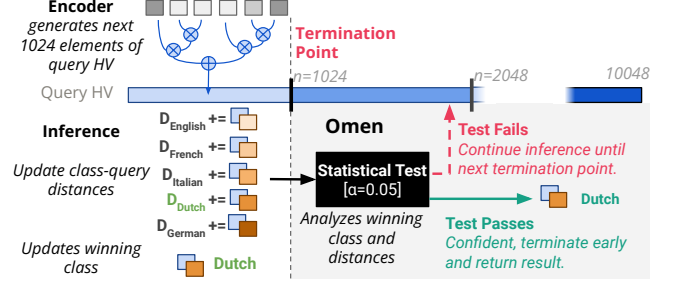
**HDC Variant.** This HDC classifier architecture is general and agnostic to the HDC variant or the hypervector size. In this example, we instantiate the HDC classifier to use the BSC variant and to work with  $N=10,048$ -element BSC hypervectors. Since each BSC element is a 0/1 value, each hypervector is 10,048 bits. We choose a hypervector size of 10,048 because it is a multiple of 64 and convenient for bit packing optimizations. The Hamming distance  $HD(V, V')$  is used to compute distances between class and query hypervectors, and the xor, majority vote, and rotational shift operations implement binding, bundling and permutation.

**Dataset.** In this example, we choose 10,000 sentences each from 5 languages, Dutch, English, French, German, and Italian, from the Leipzig Corpora Collection [21] language dataset. For each language, a random selection of 6/7 of all sentences is used as training samples, and the others as testing samples. We use only 5 languages to keep the example simple and tractable. The HDC language classifier achieves 98.5% test accuracy.

### 3.2 Challenges with Inference

As an example, we use the HDC classifier to classify two sentences Q1, Q2, where both sentences are Dutch:

- Q1. Door puur geluk leef je nog.
- Q2. Dringend burgemeesters uit hun ambt zetten als hoofd van de politie.



(a) Early termination with OMEN.

	Size	$V_{\text{Dutch}}$	$V_{\text{English}}$	$V_{\text{French}}$	$V_{\text{German}}$	$V_{\text{Italian}}$
$V_{Q1}$	10048	<b>0.4790</b>	0.5001	0.4792	0.4916	0.4831
$V_{Q2}$	10048	<b>0.4033</b>	0.4534	0.4482	0.4388	0.4690
$V_{Q1}$	9048	0.4799	0.4997	<b>0.4798</b>	0.4911	0.4809
$V_{Q2}$	9048	<b>0.4048</b>	0.4546	0.4465	0.4418	0.4719

(b) Hamming Distances of Q1, Q2 queries. Minimum distances are marked as bold (correctly classified) or red (incorrectly classified).

$p$ -values computed by the Wald’s test			sorted by $p$ -value (descending order) and compared to $\alpha/j$ (Holm-Bonferri method)		
name	$p_j$	$j$	name	$p_j < \alpha/j$	result
$H_{\text{English}}^0$	0.007560	1	$H_{\text{French}}^0$	$0.050722 \geq 0.0500$	FAIL
$H_{\text{French}}^0$	0.050722	2	$H_{\text{English}}^0$	$0.007560 < 0.0250$	PASS
$H_{\text{Italian}}^0$	0.000067	3	$H_{\text{German}}^0$	$0.004512 < 0.0166$	PASS
$H_{\text{German}}^0$	0.004512	4	$H_{\text{Italian}}^0$	$0.000067 < 0.0125$	PASS

(c) OMEN’s statistical test as applied to query Q2 for  $n = 1024$ . The unsatisfied test is marked as red.

**Figure 2.** OMEN classifier with early termination. Tables report class-query hypervector distances and statistical tests for Q1, Q2 queries.

Figure 2b presents the hypervector distances for the Q1 and Q2 classification tasks. Using the entire 10,048-element hypervectors, the Q1 and Q2 queries are both correctly classified as Dutch sentences.  $V_{Q1}$  is almost as close to  $V_{\text{French}}$  as to  $V_{\text{Dutch}}$ , and  $V_{Q2}$  is much closer to  $V_{\text{Dutch}}$  than other  $V_{\text{lang}}$ ’s.

We can reduce the memory and footprint of the model and the computation required for inference by reducing the size of the hypervectors  $N$  in the HD classifier. This optimization sacrifices classification accuracy for resource usage. If we were to repeat the above inference with the BSC language classifier, but with slightly smaller 9,048-bit hypervectors, the Q1 query would be classified incorrectly, but the Q2 query would still be classified correctly. If we investigate the hypervector distances on this reduced-size model, we find that  $V_{Q1}$  becomes closest to  $V_{\text{French}}$ , although by a very small margin, resulting in a misclassification. In contrast, the  $V_{Q2}$  is robustly much closer to  $V_{\text{Dutch}}$  and is therefore still classified correctly.

The crux of the issue is that some queries are close to multiple classes and thus difficult to classify, requiring more bits to reach a conclusion, whereas other queries are much closer



to one class and thus easy to classify. Therefore, it would be productive to dynamically adjust the number of processed bits depending on the difficulty of the individual inputs.

### 3.3 Early Termination with OMEN

OMEN enables dynamic early termination of HD classifiers for individual inputs while providing an upper bound on the accuracy loss. OMEN’s workflow is shown in Figure 2a. OMEN accepts a list of early termination points: for the current example, termination points are  $n=1024, 2048, 4096$ , and  $9192$  hypervector elements.

OMEN keeps track of the running query-class distances and tests if the current classification result is statistically significant at each termination point. For Q2, at the termination point 1 where  $n = 1024$ , it is closest to class Dutch. To test the the statistical significance of the result, OMEN applies a *statistical test* to the partially processed query and class hypervectors and the running query-class distances.

**Rationale.** OMEN frames the classification task as a *statistical inference problem*. A statistical inference problem sets hypotheses about the underlying distributions, computes  $p$ -values that quantify probabilities of observing the current samples assuming the hypotheses are true, and then infers the truthfulness of the hypotheses based on  $p$ -values.

We develop a *statistical inference* view of HDC where distances between hypervector elements are treated as samples or observations of an unknown distance distribution. The running distances between class and query hypervectors can then be interpreted as approximations of the expected value of the underlying distance distribution. We can apply statistical tests to check whether the current observations support a hypothesis about the underlying distribution. Section 4 describes this view of HDC in more detail.

**Assessing Statistical Significance.** Suppose that, at a termination point, OMEN finds that  $V_{Q2}$  is closest to  $V_{Dutch}$  (that is, a prefix of  $V_{Q2}$  is closest to the corresponding prefix of  $V_{Dutch}$ ). OMEN checks the statistical significance of the current observation. OMEN forms a hypothesis  $H_1$ :  $V_{Q2}$  is closest to  $V_{Dutch}$  in expectation. If OMEN accepts  $H_1$ , then OMEN terminates at this point. OMEN accepts a user-specified statistical significance parameter  $\alpha$  that serves as an upper bound on accuracy drop due to early termination. Therefore, if  $H_1$  is not true, the probability that OMEN accepts  $H_1$  and terminates with Dutch as the incorrect classification result should be at most  $\alpha$ . In this example, we use  $\alpha=0.05$ .

**The Statistical Test.** We next describe how the statistical test is performed. We first break  $H_1$  down into 4 hypotheses  $H_1^{English}, H_1^{French}, H_1^{German}$ , and  $H_1^{Italian}$ , which are that  $V_{Q2}$  is closer *in expectation* to  $V_{Dutch}$  than to  $V_{English}, V_{French}, V_{German}$ , or  $V_{Italian}$ , respectively. We set 4 null hypotheses that we want to reject,  $H_0^{English}, H_0^{French}, H_0^{German}$ , and  $H_0^{Italian}$ , which are negations of the previous 4 hypotheses, respectively.

The way a statistical test rejects a null hypothesis  $H_0^{lang}$  is to compute a  $p$ -value  $p_{lang}$ , which is an upper bound probability that the current distances are observed if  $H_0^{lang}$  is true. Then, if the  $p$ -value is sufficiently small (e.g.,  $p_{lang} < \alpha$ ), we can reject  $H_0^{lang}$  confidently because even if  $H_0^{lang}$  is true, we make mistakes less than  $\alpha$  of the time. OMEN uses the Wald’s test (Section 5.1) to derive the  $p$ -value for each null hypothesis. Because OMEN tries to reject all 4 null hypotheses  $H_0^{English}, H_0^{French}, H_0^{German}$ , and  $H_0^{Italian}$ , the error probability for each null hypothesis accumulates. To upper bound the overall error probability, OMEN uses the Holm-Bonferri method (Section 5.2) to derive a series of  $p$ -value thresholds for the null hypotheses from the statistical significance parameter  $\alpha$ ; briefly put, if the 4  $p$ -values from the Wald’s test are sorted into descending order, then the adjusted threshold for the  $j$ -th  $p$ -value is  $\alpha/j$ . The  $p$ -values and thresholds for query Q2 at termination point  $n=1024$  are summarized in Figure 2c.

**Continuing Execution.** The above statistical test fails because the null hypothesis  $H_0^{French}$  has a  $p$ -value 0.050722 that exceeds the Holm-Bonferri adjusted threshold 0.0500. Therefore, OMEN draws the conclusion that not enough hypervector elements have been processed to declare Dutch as the winning class and OMEN will continue processing elements until the next termination point is reached. OMEN repeats the statistical test at the second termination point after 2048 elements have been processed—at this point, the Dutch class is still the current winner. This time, the statistical test returns true, so OMEN terminate inference early and returns Dutch as the classification result. For query Q2, inference is terminated after encoding and processing 2048 hypervector elements—just 20.4% of the whole hypervector.

**OMEN is Adaptive.** OMEN’s early termination algorithm naturally adapts to differences in classification difficulty across inputs and to hardware error. For example, OMEN does not terminate early for Q1, for it is significantly harder to classify Q1 query. OMEN also naturally adjusts to any hardware errors and data corruptions. For example, if we repeat the same classification tasks but allow for bit corruptions during the distance computations, where the probability of corruption for each bit is 0.15, OMEN correctly classifies Q2 but exits 1 termination points later, after processing 40.8% of the hypervector.

## 4 Statistical View of HDC

OMEN performs early termination while also achieving a statistical upper bound  $\alpha$  on the accuracy loss. To achieve this goal, we develop an alternative statistical interpretation of HDC classification that reframes HDC distance calculations as sampling from an underlying distribution. This view of HDC classification is critical to enabling early termination with guarantees, as it enables us to apply statistical tests to evaluate the significance of a particular distance calculation. In this section, we give an intuitive overview of the statistical view and focus on the simplest HD training algorithm

Notation	Meaning
$\alpha$	significance level; upper bound of accuracy loss
$V, v$	$V$ is a hv, and $v$ is an element in $V$
$N, n$	number of total and processed hv elements $N \geq n$
$D, d$	vector of hv element distances, $d$ is an element in $D$
$P_V, P_D$	underlying distribution of each $v$ in $V$ , or $d$ in $D$
$H_1, H_0, p$	test and null hypothesis, $p$ is $p$ -value for $H_0$
$Z_n(D)$	average of first $n$ samples from $D$ , $Z(D) = \frac{1}{n} \sum_{i=1}^n D[i]$
$CD_n(D_1, D_2)$	average squared differences of $n$ -prefixes of $D_1$ and $D_2$
$\pi$	a permutation $[n] \rightarrow [n]$ for some integer $n$
$B, C$	the number of basis and class hvs, respectively
$W_i$	a vector of the $i$ -th dimensions of all $B+C$ stored hvs

**Table 2.** Formalism symbols. hv is short for hypervectors.

that makes the hypervector elements independent and identically distributed (*iid*). The discussion and formalism of the advanced training algorithms are in Section 6. We summarize our use of notations in Table 2 for reference.

#### 4.1 Fully Distributed Hypervectors

This statistical view of HDC applies if the hypervectors in the distance calculation are *fully distributed*, which intuitively means that the encoded information is fully distributed in all vector dimensions. This property of HDC was identified by prior work [37, 44]. We relate it to the statistical property *exchangeability* in Section 6. In this section, we consider a case of fully distributed hypervector, where each hypervector element  $v = V[i]$  is *iid* from the same underlying distribution  $v \sim P_V$ . We call a hypervector *iid* if its elements are *iid*. If a hypervector is *iid*, then the following intuitions hold:

- The hypervector size,  $N$ , can be thought of as the maximum number of samples that may be taken from  $P_V$ . The number of hypervector elements  $n$  that have been processed corresponds to the number of samples that have been taken so far.
- Any element-wise operations ( $\forall i, V[i] = op(V_1[i], V_2[i] \dots)$ ) over two or more independent *iid* hypervectors produces another *iid* hypervector.
- Permutation (such as rotational shifts) on *iid* hypervectors produces *iid* hypervectors.

**4.1.1 HD Operators.** All randomly generated basis hypervectors are *iid* by construction. From this point, binding, bundling, and permutation operations all produce *iid* result hypervectors. Note that certain compositions of HD operators over random hypervectors produce hypervectors that are not *iid*, but the introduced element dependence does not affect the distance properties much, and these compositions can be replaced with alternatives that do not introduce dependence while achieving the same distance properties [72]. Therefore, we may assume compositions of the HD operators over random hypervectors produce *iid* hypervectors.

**4.1.2 Element Distances/Similarities.** In HDC, hypervector distances/similarities are computed by aggregating

over element-wise distances/similarities. Note that our statistical view applies to both distance and similarity measures, but we will present them as “distance” to be concise. For hypervectors  $V_1$  and  $V_2$ , we can compute an element distance vector  $D$  where  $D[i] = dis(V_1[i], V_2[i])$ . Because the *dis* operator is element-wise,  $D$  is *iid* if  $V_1, V_2$  are *iid*. We denote each element  $d$  in  $D$  is from distribution  $P_D$  ( $d \sim P_D = dis(P_{V_1}, P_{V_2})$ ). For BSC,  $dis(V_1[i], V_2[i]) = V_1[i] \text{ xor } V_2[i]$ , the bit difference and for MAP,  $dis(V_1[i], V_2[i]) = V_1[i] V_2[i]$ , the product.

#### 4.2 Distances as Expected Value Estimates

Given a *iid* distance vector  $D$  where  $d \sim P_D$ , the HDC distance metric  $Z_n(D) = 1/n \cdot \sum_{i=1}^n D[i]$  estimates the expected element distance  $E(P_D)$  by computing the average element distance over  $n$  processed samples. Note that the expected distance and the distance distribution are both *unknown* and has no closed-form solution in complex tasks such as classification and therefore can only be estimated by sampling the distribution. The number  $n$  of hypervector elements processed corresponds to the number of samples drawn from the underlying distributions and thus affects the quality of approximation. Conceptually,  $Z_n(D) \rightarrow E(P_D)$  when  $n$  is sufficiently large.

**4.2.1 Vanilla HD Classification.** In the simplest HD classification workflow, the encoder uses HD operators on randomly generated hypervectors, and the vanilla training algorithm bundles hypervector representations of training samples from a class as the class hypervector. Thus, the class hypervectors are *iid*. During inference, the classification of a test sample is based on the distances between the test sample hypervector and the class hypervectors. Since all hypervectors involved are *iid*, the associated distance vectors are also *iid*.

#### 4.3 Distance Comparisons with Statistical Tests

We observe that with this alternate view of HDC, distance comparisons resemble sampling and testing problems, where inferential statistics, or statistical tests, can be used to assess the statistical significance of conclusions drawn from hypervector distance comparisons. Statistical tests offer guarantees that false conclusions are drawn with a probability of at most  $\alpha$ , the user-defined significance level.

**Distance Tests.** We can interpret a given distance comparison as a statistical test. When comparing two hypervector distances computed with  $n$  elements,  $Z_n(D_1) \leq Z_n(D_2)$ , we essentially compare expected values of the underlying distance distributions  $E(P_{D_1})$  and  $E(P_{D_2})$ , with  $Z_n(D_1)$  and  $Z_n(D_2)$  as their estimates. Therefore, we can formulate a *distance test* where we are testing whether the sampled hypervector distances  $D_1, D_2$  support the test hypothesis  $H_1 : E(P_{D_1}) \leq E(P_{D_2})$ . We set the *null hypothesis*  $H_0 : E(P_{D_1}) > E(P_{D_2})$ , the negation of the hypothesis  $H_1$  we want to obtain. When  $H_0$  is true, we allow the incorrect conclusion  $H_1$  to be drawn from the samples, with a probability of at most  $\alpha$ . The statistical test computes a  $p$ -value  $p$  for  $H_0$ , where  $p$  is an upper bound on the

probability that we observe  $D_1, D_2$  when  $H_0$  is true. If  $p < \alpha$ , we reject  $H_0$  and conclude that  $H_1$  holds. Otherwise, we cannot conclude  $H_1$  although we observe  $Z_n(D_1) \leq Z_n(D_2)$ , because the observations  $D_1$  and  $D_2$  are not statistically significant.

## 5 Early Termination with OMEN

In this section, we present the technical details of OMEN that enables early termination of the inference for individual inputs while providing an upper bound on the accuracy loss. OMEN adopts the statistical inference view of HDC presented in Section 4 and therefore works with classifier architectures that produce fully distributed vectors. Sections 5.1 and 5.2 present the statistical tests used by OMEN to assess the statistical significance of the classification result and introduce the algorithmic optimizations made to improve computational efficiency. Section 5.3 presents the core inference algorithm and relates the statistical significance to accuracy loss.

### 5.1 Key Statistical Tests

OMEN uses Wald's test [68] to perform distance tests and the Holm-Bonferri method to effectively draw conclusions from a set of distance tests. They are used together to assess the statistical significance of a classification result.

**5.1.1 Distance Comparisons with the Wald's test.** The Wald's test computes a  $p$ -value  $W_n(D_1, D_2)$  for a hypothesis involving the expected value of two distributions  $H_0: E(P_{D_1}) > E(P_{D_2})$ , given  $n$  observations from  $D_1$  and  $D_2$ . The  $p$ -value expression for applying the Wald's test to  $H_0$  is

$$W_n(D_1, D_2) = 1 - \Phi\left(\frac{\sqrt{n}(Z_n(D_2) - Z_n(D_1))}{\sqrt{CD_n(D_1, D_2) - (Z_n(D_2) - Z_n(D_1))^2}}\right)$$

where  $\Phi$  is the cumulative density function (CDF) of standard normal distribution, and  $n$  is the number of bits processed.

The  $CD_n$  function is cumulative squared differences between observations (distances) from two distributions, which is used to perform a *paired comparison*.  $CD_n$  is used to account for the fact that  $D_1$  and  $D_2$  may not be independent. For example, in classification, since the test input belongs to one class, closer to one class correlates to farther from other classes.

$$CD_n(D_1, D_2) = 1/n \cdot \sum_{i=1}^n (D_1[i] - D_2[i])^2$$

### 5.1.2 Optimization: Precompute $CD_n$ for Hamming Distance.

The cumulative squared difference formula can be rewritten to eliminate a hypervector if the Hamming distance metric was used to construct  $D_1$  and  $D_2$  from hypervector elements and  $D_1 = V_1 \text{ xor } V$  and  $D_2 = V_2 \text{ xor } V$ . The cumulative squared difference can be rewritten to eliminate  $V$ , as  $\forall i, V_1(i), V_2(i), V(i)$  are 0/1 bits:

$$\begin{aligned} CD_n(D_1, D_2) &= 1/n \cdot \sum_{i=1}^n (D_1[i] - D_2[i])^2 \\ &= 1/n \cdot \sum_{i=1}^n ((V_1[i] \text{ xor } V[i]) - (V_2[i] \text{ xor } V[i]))^2 \\ &= 1/n \cdot \sum_{i=1}^n (V_1[i] \text{ xor } V_2[i])^2 \\ &= 1/n \cdot \sum_{i=1}^n (V_1[i] \text{ xor } V_2[i]) \end{aligned}$$

Furthermore, if  $V_1$  and  $V_2$  are constant, such as two class hypervectors that are fixed after training, the cumulative square difference can be pre-computed for particular  $V_1$  and  $V_2$ .

### 5.2 Holm-Bonferri Method

If only one distance comparison is performed, a statistical test, such as the Wald's test, may be used to determine the significance of the result. However, when multiple distance comparisons must be made, multiple tests must be performed to find out whether we can reject all null hypotheses, where the *multiple comparison problem* arises. It is the problem of accumulating error rate of rejecting null hypotheses when there is more than one statistical test.

The Holm-Bonferroni method [28] can be used to perform multiple tests that try to reject multiple ( $m \geq 1$ ) null hypotheses  $H_0^{(1)}, H_0^{(2)}, \dots, H_0^{(m)}$ . This algorithm guarantees that the probability of falsely rejecting *any* null hypothesis  $H_0^{(i)}$  is within an upper bound of  $\alpha$ . The algorithm achieves this by using successively stricter thresholds for the null hypotheses.

**Algorithm.** The algorithm first computes the  $p$ -value for each null hypothesis and then sorts all the null hypotheses by their  $p$ -values in descending order, thus assigning to each null hypothesis a distinct integer rank  $j$ ,  $1 \leq j \leq m$ , such that, for any  $1 \leq i_1, i_2 \leq m$ , if the rank of  $H_0^{(i_1)}$  is less than the rank of  $H_0^{(i_2)}$ , then the  $p$ -value of  $H_0^{(i_1)}$  is not less than the  $p$ -value of  $H_0^{(i_2)}$ . Then the significance test used for the null hypothesis  $H_0^{(i)}$  is  $p_i < \alpha_i = \frac{\alpha}{j}$ , where  $j$  is the rank of  $H_0^{(i)}$ . If all  $p$ -value comparisons evaluate to true, we can reject all the null hypotheses. The probability that we draw this conclusion when any of the null hypothesis is true is upper bounded by  $\alpha$ . **Optimization: rewriting  $p$ -value comparison.** We can optimize the  $p$ -value comparisons for the Holm-Bonferroni method when the  $p$ -values was computed using the Wald's test. For the  $i$ -th hypothesis with a  $p$ -value  $p_i = W_n(D_1, D_2)$  at rank  $j$ , the test can be rewritten as

$$\frac{n(nZ_n(D_2) - nZ_n(D_1))^2}{n(nCD_n(D_1, D_2)) - (nZ_n(D_2) - nZ_n(D_1))^2} > \left(\Phi^{-1}\left(1 - \frac{\alpha}{j}\right)\right)^2$$

where  $\Phi^{-1}$  is the inverse of standard normal CDF. The rewrite enables several optimizations. It multiplies  $n$  to both the numerator and the denominator in the original expression in the fraction so that each normalized term  $CD_n(D_1, D_2)$ ,  $Z_n(D_2)$ , and  $Z_n(D_1)$  is multiplied by  $n$ , and no normalization is needed during computation. It also eliminates the square root operation and relocates  $\Phi$  to the right-hand side. The  $\left(\Phi^{-1}\left(1 - \frac{\alpha}{j}\right)\right)^2$  term only involves the  $\alpha$  parameter and the  $p$ -value rank  $j$ , and can be precomputed offline and stored.

### 5.3 OMEN Inference Algorithm

The OMEN inference algorithm uses the statistical view of HDC to perform early termination. The training algorithm works with full-length hypervectors of size  $N$  and produces



```

1: procedure INFERENCE(input, encoder, Vc, N, term_points,  $\alpha$ )
2:   k = len(Vc) // number of classes
3:   D = array(k, N) // element distances
4:   dist = zeros(k) // cumulative distances to classes
5:   for n in 1..N do
6:     vq = encoder(input, n) // encode the  $n$ th dimension
7:     for i in 1..k do
8:       D[i][n] = dis(vq, Vc[i][n])
9:       dist[i] += D[i][n]
10:    cand = argmin(dist) // argmax if dis is a similarity metric
11:    if n in term_points and confident(D, cand, n,  $\alpha$ ) then
12:      return cand
13:   return cand

```

**Figure 3.** OMEN’s inference algorithm.  $N$  is hypervector size.  $V_c$  is class hypervectors.  $\alpha$  is confidence level.  $cand$  is current inference result. **dis** returns distance between two hypervector elements. **confident** uses statistical tests.

a set of fully distributed class hypervectors. OMEN deploys an HDC classifier that works with a fully distributed query hypervector  $V_q$  that is computed from the provided HD encoder.

**5.3.1 Algorithm.** OMEN’s algorithm is inspired from Wald’s sequential probability ratio test (SPRT), which sequentially samples and terminates an experiment when the evidence is statistically significant in manufacturing quality control [68]. The SPRT is used under manufacturing conditions where sampling products is costly. We draw a connection between computational cost and product sampling cost.

Figure 3 shows the pseudo code for OMEN’s inference algorithm. OMEN accepts as input a set of termination points, where each termination point  $j$  is reached after computing over  $n_j$  hypervector elements. We denote the number of termination points as  $m$ . OMEN breaks up the inference task into  $m + 1$  units of work and checks if the inference task can be terminated early after each of the first  $m$  units of work. After a given work unit has been completed at the  $j$ -th termination point,  $n_j$  elements of  $V_q$  have been generated by the encoder. OMEN computes the query-class element distances. After each unit of work is completed, OMEN finds the current closest class **cand** to the query and assesses its statistical significance.

Next, OMEN assesses the statistical significance of this classification result using the Holm-Bonferroni method from Section 5.2 using the user-provided error threshold  $\alpha$ . If the Holm-Bonferroni method returns **true**, OMEN terminates inference early and returns class **cand**. If the method returns **false**, OMEN continues execution. If all statistical tests fail, OMEN eventually returns the classification result after all  $N$  dimensions have been processed. OMEN can optionally be configured to mark inputs that fail all statistical tests as difficult samples and pass them on to more reliable and sophisticated classifiers.

**5.3.2 Guarantees on Classification Accuracy Loss.** The statistical level  $\alpha$  is an upper bound on the drop in classification accuracy resulting from processing fewer than  $N$  hypervector elements. Misclassifications that occur even with

infinite hypervector size are attributed to *model error*. Specifically, *model error* occurs when  $\exists i \neq t, E(P_{D_i}) < E(P_{D_t})$ , where  $D_i$  is distance to the  $i$ -th class and  $t$  is the ground-truth label. Intuitively, this means the query is expected to be closer to other class than the ground-truth class. This form of error can only be mitigated with better encoder and classifier design.

We denote the model error rate *mer* as the percentage of test inputs in the test distribution that suffers from model error. Misclassifications that occur when the hypervector size is not large enough are attributed to *estimation error*. This happens when  $\forall i \neq t, E(P_{D_i}) \geq E(P_{D_t})$  but the returned classification  $cand \neq t$ . Because  $E(P_{D_{cand}}) \geq E(P_{D_t})$  but we get  $Z_n(D_{cand}) < Z_n(D_t)$ , this kind of error can be mitigated by computing more distance samples. The *estimation error rate* is therefore dependent on the number  $n$  of dimensions processed, which we denote as *eer*( $n$ ). The expected accuracy, when always processing all  $N$  hypervector dimensions without early termination, can therefore be defined as

$$acc(N) = 1 - mer - eer(N)$$

OMEN may terminate early, and may suffer from larger estimation error rate. Denoting  $err_{omen}(N, \alpha)$  as the estimation error rate OMEN has with confidence level  $\alpha$  and hypervector size  $N$ , we can define OMEN’s accuracy as

$$acc_{omen}(N, \alpha) = 1 - mer - err_{omen}(N, \alpha).$$

We also define  $acc_{opt} = \lim_{N \rightarrow \infty} acc = 1 - mer$ , which is the best accuracy we can get out of the current model if we were to use infinitely large hypervectors and be free of estimation error. Because we have  $err_{omen}(N, \alpha) > err(N) > 0$ , the accuracy relationship between the early-termination, normal inference, and optimal inference is  $acc_{omen}(N, \alpha) < acc(N) < acc_{opt}$ . The Holm-Bonferroni method guarantees that when  $\forall i \neq t, E(P_{D_i}) \geq E(P_{D_t})$  (i.e., the input does not suffer from model error), we draw conclusion of  $\forall i \neq cand, E(P_{D_i}) \geq E(P_{D_{cand}})$  for some  $cand \neq t$  with a probability at most  $\alpha$ . This translates to that  $\alpha$  is an upper bound on OMEN’s estimation error rate, i.e.,  $err_{omen}(N, \alpha) \leq \alpha$ . Therefore, we have:

$$\begin{aligned}
acc_{omen}(N, \alpha) &= 1 - mer - err_{omen}(N, \alpha) \\
&\geq 1 - mer - \alpha \\
&\geq acc_{opt} - \alpha \\
&\geq acc(N) - \alpha
\end{aligned}$$

In summary, the statistical test employed by OMEN ensures that the accuracy loss due to early termination over using the whole vector is at most  $\alpha$ , or  $acc(N) - acc_{omen}(N, \alpha) \leq \alpha$ .

## 6 Exchangeability

While Section 4 intuitively explains the *iid* property of HDC in the simple case, the state-of-the-art training algorithms often break the independence between the hypervector elements as they use aggregate information of all elements (e.g., hypervector distances) to update them. We find that although the learned hypervectors are not *iid*, they are *exchangeable*, which is a weaker statistical property than *iid* and relates nicely to the *fully distributed* property. In this section, we formalize exchangeability in Section 6.1, prove that the two



representative training algorithms produce exchangeable hypervectors in Section 6.2 and discuss the relation to the *fully distributed* property and Wald’s test in Section 6.3.

### 6.1 Definition of Exchangeability

We formally define *exchangeability* in this section. We first define *permutation* and *symmetric functions* as follows.

**Definition 1** (Permutation). Denote  $m \in \mathbb{Z}$  as an integer, and  $[m]$  as the set of integers from 1 to  $m$ . A permutation  $\pi$  is a function  $[m] \rightarrow [m]$ , where  $\{\pi(i) : i \in [m]\} = [m]$ . Because  $\pi$  is a one-to-one mapping from  $[m]$  to  $[m]$ , it is invertible, and we denote the inverse of  $\pi$  as  $\pi^{-1}$ . For a vector of  $n$  elements  $W = (W_1, W_2, \dots, W_n)$ , we denote  $W$  shuffled by  $\pi$  as  $\pi W = (W_{\pi(1)}, W_{\pi(2)}, \dots, W_{\pi(n)})$ .

**Definition 2** (Symmetric Function). A function  $f(X)$ , where  $X = (x_1, \dots, x_n)$  is a  $n$ -element vector, is symmetric about  $X$  if  $\forall \pi, f(\pi X) = f(X)$ .

The definition of exchangeability is taken from Kuchibhotla [46] (real-valued random variables, assuming the probability density function exists).

**Definition 3** (Exchangeability). Suppose  $W = (W_1, \dots, W_n) \in \mathcal{R}^n$  is a vector of real-valued random variables,  $W_i$ ’s are exchangeable iff their joint probability density function  $p(W)$  is symmetric about  $W$ .

iid implies exchangeable because if  $W_i$ ’s are iid, we have  $p(W) = \prod_i p(W_i) = \prod_i f(W_i)$ , where  $f$  is the identical probability density function shared by all the  $W_i$ ’s.  $p(W)$  function is symmetric about  $W$  because  $p(\pi W) = \prod_i f(W_{\pi(i)}) = \prod_i f(W_i) = p(W)$  due to the commutativity of multiplication. However, the inverse is generally not true, because there can be dependence among  $W_i$ ’s.

### 6.2 Learned Hypervectors Are Exchangeable

In this section, we formalize some representative advanced HDC training algorithms and prove that they produce exchangeable hypervectors.

**6.2.1 Primitives.** We first define some symmetry-related concepts, and prove some useful theorems.

**Definition 4** (Element-wise Function). A function  $f(X) = Y$  where  $X, Y$  are  $n$ -element vectors is called *element-wise* if  $f(\pi X) = \pi f(X)$ .

**Lemma 1.** Composition of element-wise functions produces element-wise functions.

*Proof.* If  $f$  and  $g$  are element-wise, for any  $n$ -element vector  $X$ , we have  $f(g(\pi X)) = f(\pi g(X)) = \pi f(g(X))$ , which implies that  $f \circ g$  is element-wise.  $\square$

**Lemma 2.** Composition of symmetric and element-wise functions are symmetric.

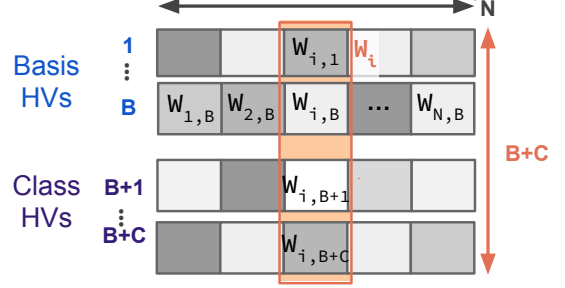


Figure 4. HD representation -  $W$ ’s definition.

*Proof.* If  $f$  is symmetric and  $g$  is element-wise,  $f \circ g$  is symmetric because  $\forall X, f(g(\pi X)) = f(\pi g(X)) = f(g(X))$ , and  $g \circ f$  is also symmetric because  $\forall X, g(f(\pi X)) = g(f(X))$ .  $\square$

**Lemma 3.** For a function  $f(X)$  symmetric about  $X$ ,  $f$ ’s gradient  $\nabla f(X) = (\frac{\partial f}{\partial x_1}(X), \dots, \frac{\partial f}{\partial x_n}(X))$  satisfies  $\nabla f(\pi X) = \pi \nabla f(X)$ .

*Proof.*  $\nabla f(\pi X) = (\frac{\partial f}{\partial x_1}(\pi X), \dots, \frac{\partial f}{\partial x_n}(\pi X))$ . We have  $\forall i \in [n]$ ,

$$\begin{aligned} \frac{\partial f}{\partial x_i}(\pi X) &= \lim_{\epsilon \rightarrow 0} 1/\epsilon \cdot (f(\pi X + \epsilon e_i) - f(\pi X)) \\ &= \lim_{\epsilon \rightarrow 0} 1/\epsilon \cdot (f(X + \epsilon e_{\pi(i)}) - f(X)) \\ &= \frac{\partial f}{\partial x_{\pi(i)}}(X) \end{aligned}$$

where  $e_i$  is a  $n$ -element vector with the  $i$ -th element being 1 and others being 0. The first and third equalities are the definition of the derivative, and the second equality exploits the symmetry of  $f(x)$ . Note that  $\pi^{-1}e_i$  is  $e_{\pi(i)}$  instead of  $e_{\pi^{-1}(i)}$  (Definition 1). This is equivalent to  $\nabla f(\pi X) = \pi \nabla f(X)$ .  $\square$

Because the hypervectors are initialized to be exchangeable (e.g., zeros or random Gaussian), we only need to prove that each update of the training algorithms *preserves exchangeability*. A transformation from a vector of random variables to another vector of random variables is *exchangeability preserving* iff for all exchangeable input random variables, the transformed output random variables are also exchangeable. The following theorem states the condition for which a transformation on a vector of exchangeable random variables preserves exchangeability [11, 46].

**Theorem 1** (Exchangeability Preservation). Suppose  $W = (W_1, \dots, W_n) \in \mathcal{W}^n$  is a vector of exchangeable random variables. Fix a transformation  $G: \mathcal{W}^n \rightarrow (\mathcal{W}')^m$ . If for each permutation  $\pi_1: [m] \rightarrow [m]$  there exists a permutation  $\pi_2: [n] \rightarrow [n]$  such that

$$\pi_1 G(W) = G(\pi_2 W), \quad \text{for all } W \in \mathcal{W}^n$$

then  $G(\cdot)$  preserves exchangeability of  $W$ .

**Corollary 1.** Element-wise functions are exchangeability preserving transformations.

*Proof.* If  $f: \mathcal{W}^n \rightarrow (\mathcal{W}')^n$  is element-wise, for each  $\pi_1: [n] \rightarrow [n]$ , let  $\pi_2 = \pi_1$ , then we have  $\pi_1 f(W) = f(\pi_2 W)$ ,  $\forall W \in \mathcal{W}^n$ . By Theorem (1),  $f$  preserves exchangeability.  $\square$

**6.2.2 HDC Representation  $W$ .** Our insight is that although the HDC training algorithms update hypervectors in a complex way, they do not break the symmetry among hypervector elements for all the hypervectors. To show that, we need to consider all hypervectors we have (basis hypervectors in codebooks and class hypervectors) as a whole. Denote  $N$  as our hypervector length,  $W_i$  as a vector of the  $i$ -th dimensions of all  $B$  basis hypervectors in codebooks and  $C$  class hypervectors. We illustrate  $W$ 's definition in Figure 4. Using the same notation as theorem (1), denote  $W = (W_1, \dots, W_N)$ .  $W$  is initialized iid. We demonstrate that the training algorithms are exchangeability preserving transformations on  $W$ .

**6.2.3 Implications of  $W$ 's Exchangeability.** In this section, we prove that when the HDC representation  $W$  defined above is exchangeable, the encoding produces exchangeable hypervectors, and their distances to the class hypervectors are symmetric about  $W$ .

**Theorem 2** (Exchangeable Encoding). *When  $W$  is exchangeable, HD encoding  $enc(W, T)$  produces a vector of exchangeable elements, where  $T$  is an input data to encode.*

*Proof.* HD encoding is a function  $enc(W, T)$  of  $W$  (specifically only the codebook part  $W_{j,i}, i \leq B$ ) and an input  $T$ .  $enc$  is a composition of HD operators binding, bundling, and permuting. They are element-wise functions on the codebook.<sup>1</sup> Therefore, by Lemma (1),  $enc$  is an element-wise function. By Corollary (1),  $enc$  is exchangeability preserving, so when  $W$  is exchangeable, its output hypervector is also exchangeable.  $\square$

**Corollary 2** (Exchangeable Element Distances). *When  $W$  is exchangeable, element distances  $ed(W, T)$  between an encoded hypervector  $enc(W, T)$  and class hypervectors are exchangeable.*

*Proof.* Denote  $ed(W, T)_{j,i} = dis(enc(W, T)_j, W_{j,B+i})$  as the  $j$ -th element distance from the encoded hypervector to the  $i$ -th class hypervector. Since  $enc(W, T)$  is element-wise, we have

$$\begin{aligned} ed(\pi W, T)_{j,i} &= dis(enc(\pi W, T)_j, (\pi W)_{j,B+i}) \\ &= dis((\pi enc(W, T))_j, (\pi W)_{j,B+i}) \\ &= (\pi ed(W, T))_{j,i} \end{aligned}$$

Therefore,  $ed(W, T)$  is by definition also an element-wise function (Lemma (1)) and preserves exchangeability (Corollary (1)), and the result length- $N$  vector is exchangeable when  $W$  is exchangeable.  $\square$

**Corollary 3** (Symmetric Distance). *When  $W$  is exchangeable, distances  $d(W, T)$  between an encoded hypervector and the class hypervectors are symmetric about  $W$ .*

*Proof.* Denote  $d(W, T)$ , a vector of length  $C$ , as the HD distances from a encoded hypervector  $enc(W, T)$  to the  $C$  class hypervectors, that is,  $d(W, T)_i = 1/N \cdot \sum_j ed(W, T)_{j,i}$ .  $d(W, T)$

<sup>1</sup>Strictly speaking, permuting is not an element-wise function by our definition. However, because it is a special kind of exchangeability preserving function, has similar properties as element-wise functions, and can be replaced by binding in the HD encoding [44], we omit this detail to keep the proof simple.

a composition of the element-wise function  $ed(W, T)$  and the symmetric function of averaging. Therefore, by Lemma (2),  $d(W, T)$  is a symmetric function about  $W$ .  $\square$

**6.2.4 Training Preserves Exchangeability.** We now formalize a distance-based iterative training algorithm and a gradient descent-based training algorithm, and prove that they preserve the exchangeability of  $W$ .

**Distance Based.** In a distance-based iterative training algorithm, each transformation takes a training sample  $T$  with label  $t$ , encodes  $T$  into a hypervector  $enc(W, T)$  using the codebook, adds it to the class hypervectors with a weight vector  $wt(W, T, t)$  of length  $C$ . The weight applied to the hypervector of class  $t'$  is  $wt(W, T, t)_{t'}$ , which is computed from the distance  $d(W, T)_{t'}$  and whether  $t' = t$ . Usually,  $wt(W, T, t)_{t'} \geq 0$  when  $t' = t$ , and  $wt(W, T, t)_{t'} \leq 0$  when  $t' \neq t$ . OnlineHD training algorithm [26] is an example of this kind. It selects the update weights to prevent model saturation. The transformations can be formalized as follows. From the way the weights are computed, they can be rewritten as  $wt(W, T, t) = \sigma(d(W, T), t)$ , where  $\sigma$  is a function that depends on the concrete update strategies and hyperparameters like learning rate. Because  $d(W, T)$  is symmetric about  $W$  by Corollary (3),  $wt(\pi W, T, t) = \sigma(d(\pi W, T), t) = \sigma(d(W, T), t) = wt(W, T, t)$ . Therefore,  $wt(W, T, t)$  is also symmetric about  $W$ . The transformation  $G$  adds the encoded hypervector scaled by the weights to the class hypervectors

$$G(W)_{j,B+i} = W_{j,B+i} + wt_i(W, T, t) \cdot enc(W, T)_j.$$

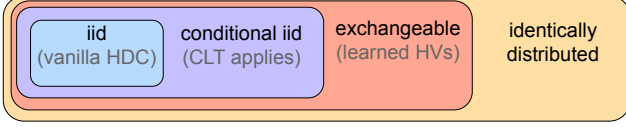
Now we prove that  $G$  is element-wise. For each  $\pi: [N] \rightarrow [N]$ ,

$$\begin{aligned} (\pi G(W))_{j,B+i} &= G(W)_{\pi(j),B+i} \\ &= W_{\pi(j),B+i} + wt_i(W, T, t) \cdot enc(W, T)_{\pi(j)} \\ &= (\pi W)_{j,B+i} + wt_i(\pi W, T, t) \cdot enc(\pi W, T)_j \\ &= G(\pi W)_{j,B+i}. \end{aligned}$$

The second to last equation makes use of the fact that  $wt_i$  is a symmetric function about  $W$  and  $enc$  is an element-wise function. For  $i \leq B$ , since  $G$  does not change the  $W_{j,i}, i \leq B$  (codebook), we also have  $(\pi G(W))_{j,i} = (\pi W)_{j,i} = G(\pi W)_{j,i}$ . By Corollary (1),  $G$  preserves exchangeability.

**Gradient Descent.** In a gradient descent-based training algorithm, each transformation uses the gradient of a loss function to update the hypervectors. The loss function is of the form  $loss(W, T, t) = \sigma'(d(W, T), t)$ , where  $T$  is a training input with label  $T$  and  $\sigma'$  usually computes the loss function (e.g., cross-entropy) comparing the distance vector  $d(W, T)$  applied by softmax to the ground-truth one-hot vector  $e_t$ . Due to the same reason as  $wt(W, T, t)$ ,  $loss(W, T, t)$  is also symmetric about  $W$ . There are different variants of gradient descent-based training algorithms. LeHDC [14] freezes the basis hypervectors and learns class hypervector representations. LDC [15] learns the basis hypervectors together with the class hypervectors. We prove a general theorem for them.

**Theorem 3.** *When the loss function is symmetric about  $W$ , gradient descent updates on  $W$  are exchangeability-preserving.*



**Figure 5.** Relationship between related statistical concepts. Central limit theorem (CLT) applies to conditional *iid* variables, which can approximate exchangeable variables.

*Proof.* Because  $\text{loss}(W, T, t)$  is also symmetric about  $W$ , by Lemma (3), we have  $\nabla_W \text{loss}(\pi W, T, t) = \pi \nabla_W \text{loss}(W, T, t)$  for any  $\pi : [N] \rightarrow [N]$ . Therefore, denoting  $\text{grad}_{j,i}(W, T, t) = \frac{\partial \text{loss}(W, T, t)}{\partial W_{ji}}$  as the gradient for  $W_j$  at its  $i$ -th element, we have  $\text{grad}_{(\pi(j)),i}(W, T, t) = \text{grad}_{j,i}(\pi W, T, t)$  for any  $\pi$ . Then, each update in the gradient descent-based training algorithm is a transformation  $G$  as follows

$$G(W)_{j,i} = W_{j,i} - lr \cdot \text{grad}_{j,i}(W, T, t)$$

where  $lr$  is a constant learning rate. We prove that  $G$  is element-wise. For each  $\pi : [N] \rightarrow [N]$ , we have

$$\begin{aligned} (\pi G(W))_{j,i} &= G(W)_{\pi(j),i} \\ &= W_{\pi(j),i} - lr \cdot \text{grad}_{\pi(j),i}(W, T, t) \\ &= (\pi W)_{j,i} - lr \cdot \text{grad}_{j,i}(\pi W, T, t) \\ &= G(\pi W)_{j,i}. \end{aligned}$$

Therefore, by Corollary (1),  $G$  preserves exchangeability.  $\square$

The theorem naturally applies to LDC. For LeHDC, since it only updates the class hypervectors ( $W_{j,B+i}, 1 \leq i \leq C$ ), we set  $\text{grad}_{j,i}(W, T, t)$  to 0 for  $1 \leq i \leq B$ , which does not break Lemma (3). The proof works with this modification.

### 6.3 Discussion

In this section, we discuss the relation between exchangeability and the fully distributed property, and the applicability of Wald’s test to exchangeable hypervectors. In Figure 5, we show a Venn diagram that illustrates the relationship between exchangeability and related concepts.

**6.3.1 Relation to Fully Distributed Property.** The *fully distributed* property of HDC intuitively means that the information encoded by the hypervector is equally distributed among the hypervector elements. Each element is not more or less important than or different from the other elements, while all elements as a whole represent the information. One may think that dependence-introducing training algorithms would break the fully distributed property and thus make the model less error-resilient. Counter-intuitively, prior work found that the training algorithms using gradient descent do not make the HDC classification model less error resilient [15]. This finding indicates that, while the hypervectors generated by these correlation-introducing training algorithms are not *iid*, the information is still fully distributed in them. Therefore, *iid* is possibly too strong and not a necessary condition for the fully distributed property. The question then arises, is there

a less restrictive statistical property that satisfies the fully distributed property that makes the model error resilient?

We postulate that the *exchangeability* property offers the *fully distributed* property without the strict independence constraint. Exchangeable hypervectors have identically distributed elements. We believe the *identically distributed* property is necessary to ensure the hypervector is fully distributed because it intuitively specifies that each element is sampled from the same information distribution. Non-identically distributed hypervectors may have elements sampled from different distributions that contain disjoint information, and not all information is guaranteed to have the necessary redundancy that makes the representation error resilient.

However, we believe that the identically distributed property alone is not sufficient to ensure the vector is fully distributed, as it does not guarantee the information is *evenly* distributed throughout the vector. Consider an example where we have a hypervector  $V$ , and each element  $v \sim P_V$  is identically distributed. If  $V[1]$  and  $V[2]$  are independently drawn from  $P_V$ , while other  $V[i]$ ’s for  $2 < i \leq n$  are copies of  $V[2]$ , then  $V[1]$  is more important and contains as much information about the underlying distribution as all the other elements combined. In fact, the information entropy of  $V[1]$  is equal to that of all other  $V[i]$ ’s, implying information is not evenly distributed throughout the vector. The exchangeability property ensures the elements in the hypervector are symmetric, since elements may be permuted without changing the joint distribution function.

**6.3.2 Assumptions Made by Wald’s Test.** We next provide a theoretical justification for applying Wald’s test to exchangeable partial distances in HD inference. The Wald’s test assumes that the samples follow the central limit theorem (CLT), such that the test statistic converges to a normal distribution. The hypervectors and distances described in Section 6.2 are exchangeable; however the CLT does not generally hold for exchangeable random variables.

We describe two methods to show CLT holds for exchangeable random variables. First, exchangeable random variables follow the CLT if the covariance of each pair of variables is zero, as this implies the variables are conditionally *iid* [7, 39]. Second, finite exchangeable random variables can be approximated with conditionally *iid* random variables [13], which follow CLT. We note the second method is always applicable, but the approximation bound is loose. We leave examining whether the first method applies to learned hypervectors to future work, this method is desirable as it does not introduce an approximation bound. We stress that, empirically, we find the Wald’s test effectively and robustly bounds the accuracy loss by  $\alpha$  (Section 7) for learned classifiers.

## 7 Evaluation

We evaluate OMEN’s accuracy and performance on 19 benchmarks and compare with the unoptimized baseline and 3 other



dataset	(train,test,class)	description	NN acc
lang	(42855,7145,5)	language classification [21]	N/A
ucihar	(7352,2947,6)	human activity recognition [3]	94.02
isolet	(6237,1560,26)	voice recognition [4]	95.50
mnist	(60000,10000,10)	image classification [50]	95.83

**Table 3.** Dataset configurations. ucihar, mnist, isolet neural-network accuracy is from MI-ShaRNN [12], Finn [67], and a DNN baseline evaluated by Cano et al. [26].

Algorithm	Encoder	Class HV
OnlineHD (OHD) [26]	random	weighted bundle
LeHDC (LeH) [14]	random	learned
LDC [15]	learned	learned

**Table 4.** Hypervector (HV) Training algorithm descriptions.

early termination heuristics. We show the detailed experimental setup in Section 7.1, the main accuracy and performance comparisons to the unoptimized baseline in Section 7.2, and comparisons to heuristics in Section 7.3. We also explore the error resiliency of our baselines in Section 7.4.

## 7.1 Experimental Setup

**Datasets.** Following prior work on HDC [15, 26], we evaluate on four datasets fit for edge-inference scenarios. The dataset configurations are summarized in Table 3. We report the number of training and testing samples and the number of classes. We also report the accuracy achieved by lightweight neural network-based models on each dataset except for language. The language dataset is described in Section 3 and was not used in the same way in prior work. We use the trigram encoder described in Section 2.1 for the language dataset. For the other datasets, we use positional encoding, the hypervector encoding methods applied in previous work [26, 60].

**HDC Training Algorithms.** To demonstrate the applicability of OMEN to various training algorithms, we evaluate on three algorithms on both BSC and MAP. The algorithms are summarized in Table 4 and are described in Section 6.2.4. Each algorithm only handled one type of HDC (BSC/MAP), and we extend them to the other type for our evaluation. Generally, learning-based training algorithms can achieve higher accuracy with fewer dimensions. To account for this difference among training algorithms and mimic real deployment scenarios where a reasonable accuracy should be achieved, we use hypervector length  $N = 10048$  for OHD,  $N = 5056$  for LeH, and  $N = 256$  for LDC. We select lengths divisible by 64 to enable bit-packing optimizations for BSC. We use a fixed random seed value 0 at the beginning of the experiment script.

**Baselines.** We summarize our 5 baselines in Table 5. For all the baselines, we use the same set of termination points. The termination points are arithmetic sequences parameterized by an initial term  $it$  and a common difference term  $ct$ , featuring checking termination condition after every  $ct$  dimension is processed, starting from dimension  $it$ . For low-dimensional

Baseline	Termination Strategy
UNOPTIMIZED	after all dimensions are processed
OMEN	the results are statistically significant
SV	process same # dimensions as average OMEN $\alpha = 5\%$
DIFF [9]	winner class' similarity margin exceeds a threshold
ABSOLUTE [30]	similarity exceeds a threshold
MEAN [53]	similarity exceeds a winner-dependent threshold

**Table 5.** OMEN and baseline termination strategies.

LDC training algorithms, we use  $it = 16$  and  $ct = 4$ . For high-dimensional LeH and OHD, we use  $it = 512$ ,  $ct = 64$  for BSC and  $it = 128$ ,  $ct = 64$  for MAP. We select  $it$  to have enough samples for valid statistical tests. We evaluate OMEN with 3 confidence levels  $\alpha = 1\%, 5\%$  or  $10\%$ , featuring conservative, balanced, and aggressive optimization modes. To compare OMEN's termination strategy with naively using fewer dimensions, we include a baseline SV (Smaller Vector) that only uses  $n$  dimensions per inference, where  $n$  is the average number of dimensions used by OMEN  $\alpha = 5\%$  on the same benchmark. For the SV baseline, we only measure its accuracy.

We also compare OMEN against heuristic early termination methods. DIFF and ABSOLUTE use one user-provided threshold respectively, which dictates the trade-off between accuracy and performance. Since these thresholds are dataset-dependent, we tune a reasonable threshold for each dataset using binary search. The tuned thresholds perform aggressive early termination while allowing a maximum 1% accuracy loss in training data. Morris et al. proposed to derive a threshold for each class from training data [53]. They computed class-dependent thresholds based on the distance statistics of correctly classified training samples and the distance statistics of misclassified training samples. Their threshold formulas are not suitable for our evaluation because in some cases, we find no misclassified training samples with a certain label. We adapt their threshold formula to obtain the MEAN baseline. MEAN uses a threshold table derived from training data. Specifically, for each class  $t$ , it computes the average of (1)  $mean + std$  of distances from training samples with label  $t$  to the hypervector of class  $t$ , and (2)  $mean - std$  of distances from training samples with label  $t' \neq t$  to the hypervector of class  $t$ , as the absolute similarity threshold when class  $t$  is the current winner.

**Execution Setup.** We implement the baselines in C++ to measure the inference performance. For OMEN, we precompute the squared distance differences  $CD$  for BSC at set termination points (described in Section 5.1.2), and lazily evaluate  $CD$  for MAP. We also apply bit-packing and precomputation optimizations to the heuristic baselines. Specifically, we precompute class statistics required for the normalized similarity metric used in the DIFF, ABSOLUTE and MEAN baselines. We conduct the evaluation on a Microcontroller (MCU) for the models that can fit in to mimic real edge inference scenarios, and on our local desktop otherwise. We use B-U585I-IOT02A, an evaluation board with a STM32U585AI MCU. This MCU has one Arm Cortex-M33 core, 2MB of flash memory, and

786KB of SRAM. In the MCU, the flash memory and SRAM sizes limit the hypervector’s size and data types. We are able to fit all models in the MCU except the 10048-dimensional MAP OHD and the 5056-dimensional MAP LeH. The desktop machine has 64GB system memory and 24-core CPU (32 threads, Intel(R) Core(R) CPU i9-13900K with 8 performance cores @ 3.00 GHz and 16 efficiency cores @ 2.20 GHz). Our program was pinned on core 0 and is single-threaded.

**Metrics.** The absolute accuracy (acc), average number of elements processed (# dim), and average runtime in milliseconds (ms) is reported for the unoptimized HD classifier for each benchmark. For each other baseline, the relative percent accuracy loss (*AL*), reduction in number of elements processed (*DRR*), and the speedup (*SUR*) relative to the unoptimized baseline are reported. Lower numbers are better for *AL*, and higher numbers are better for *DRR* and *SUR*.

## 7.2 Comparison to Unoptimized Inference

We summarize the experiment results in Table 6. We evaluate all 5 baselines on 19 benchmarks. For language dataset, we only use OHD-BSC training algorithm because this dataset is intended as an example.

**7.2.1 Early Termination Results.** OMEN can drastically reduce the number of elements processed. Compared to the UNOPTIMIZED baseline, OMEN achieves  $1.10\text{--}5.51\times$  dimension reduction for  $\alpha = 1\%$ ,  $1.28\text{--}7.86\times$  dimension reduction for  $\alpha = 5\%$ , and  $1.45\text{--}9.88\times$  dimension reduction for  $\alpha = 10\%$ .

**7.2.2 Runtime Results.** OMEN’s dimension reduction directly translates to inference runtime improvements in most benchmarks. OMEN offers  $1.08\text{--}7.21\times$  inference speed-up for  $\alpha = 1\%$ ,  $1.24\text{--}10.04\times$  inference speed-up for  $\alpha = 5\%$ , and  $1.41\text{--}12.18\times$  inference speed-up for  $\alpha = 10\%$  on 16 out of 19 benchmarks, excluding 3 LDC-BSC benchmarks. On these 16 benchmarks, the runtime speed-up ratio *SUR* and dimension reduction ratio *DRR* can be linearly fitted with  $SUR = 1.135 \cdot DRR - 0.357$ , with a small mean squared error of 0.244, indicating that OMEN’s runtime overhead is small. We show an example average latency breakdown of the UNOPTIMIZED and OMEN  $\alpha = 5\%$  baselines in the OHD-BSC-ucihar benchmark in Table 7. In this benchmark, encoding contributes the main latency, accounting for  $> 99\%$  in both baselines. OMEN’s statistical tests account for only 0.76% of the total latency while significantly reducing the encoding overhead by 68.1% compared to the UNOPTIMIZED baseline.

We note that although OMEN is able to deliver  $1.36\times$ ,  $1.16\times$  and  $1.06\times$  speed-up with  $\alpha = 10\%$  for 3 LDC-BSC benchmarks, OMEN’s dimension reduction does not directly translate to speed-up in these benchmarks. The reason is twofold. First, the LDC-BSC benchmarks employ a small hypervector size  $N = 256$  and use optimization methods to learn very compressed hypervector representations, which leaves less room

for runtime optimizations and makes it more difficult to terminate early without losing accuracy. We can see that other baselines also fail to optimize in these 3 benchmarks, as they either do not offer any dimension reduction and make inference slower, or suffer from a huge  $> 30\%$  accuracy loss. Second, the LDC-BSC benchmarks use binary and integer operators for encoding and distance calculation, while OMEN’s statistical tests use floating-point operations (FLOPs) which can be much slower than binary operations. We note that this issue can be potentially mitigated by approximating the floating-point value comparisons with integer operations, trading the accuracy guarantee off for better performance. We leave more unsound optimizations of OMEN to future work.

**7.2.3 Accuracy Results.** OMEN robustly achieves accuracies above the accuracy lower bound (UNOPTIMIZED accuracy  $-\alpha$ ). For  $\alpha = 1\%$ , OMEN’s accuracy loss compared to the unoptimized baseline is 0.0–0.6%, with no accuracy loss on 13 benchmarks. For  $\alpha = 5\%$ , OMEN’s accuracy loss is 0.0–2.4%. For  $\alpha = 10\%$ , OMEN’s accuracy loss is 0.0–4.7% on all benchmarks except OHD-MAP-ucihar. In OHD-MAP-ucihar the accuracy loss is 9.2%, close to the upper bound 10%, which shows that the accuracy loss bound of OMEN can be close to tight in certain scenarios. OMEN  $\alpha = 5\%$  also achieves 1.0–7.7% less accuracy loss than SV on all benchmarks except OHD-MAP-ucihar, processing the same average number of dimensions. OMEN  $\alpha = 5\%$ ’s accuracy loss is 0.5% more than SV in OHD-MAP-ucihar due to randomness. We also note that HDC baselines achieve better accuracy than lightweight neural network-based models reported in Table 3.

## 7.3 Comparison to Heuristic Early Termination

In this section, we compare the results of OMEN with other early termination baselines DIFF, ABSOLUTE, and MEAN in Table 6. We find that OMEN consistently outperforms the other baselines in all benchmarks, usually by a large margin. We first exclude all baseline results that suffer from a  $> 10\%$  accuracy loss, as too large an accuracy drop undermines the usefulness of the classification results. Considering other baseline results, we are able to identify at least one OMEN baseline in each benchmark that is competitive in all 3 measures compared to other baselines, i.e., OMEN achieves less or comparable ( $< 0.1\%$  difference) accuracy loss, higher dimension reduction and speed-up. We bold the results of a competitive OMEN baseline in each benchmark in Table 6. Comparing the competitive OMEN baseline with other baselines in all benchmarks, OMEN delivers an additional  $0.05\text{--}6.15\times$  dimension reduction and  $0.04\text{--}5.85\times$  speed-up over the unoptimized baseline, while achieving higher or comparable accuracy. We note that among 3 early termination heuristics, DIFF is relatively stable in accuracy, with 0.0–5.0% accuracy loss, while ABSOLUTE is unstable and has  $> 30\%$  accuracy loss in 3 benchmarks, and MEAN has  $< 10\%$  accuracy loss in only 1 benchmark. We conjecture that the poor accuracy of MEAN

Benchmark	UNOPTIMIZED			OMEN $\alpha=1\%$			OMEN $\alpha=5\%$			SV	OMEN $\alpha=10\%$			DIFF			ABSOLUTE			MEAN		
	acc	# dim	time	AL	DRR	SUR	AL	DRR	SUR		AL	DRR	SUR	AL	DRR	SUR	AL	DRR	SUR	AL	DRR	SUR
OHD-B-lang	98.5	10048	62.84	0.0	5.51	7.21	0.2	7.86	10.04	3.3	<b>0.7</b>	<b>9.89</b>	<b>12.18</b>	1.0	6.65	8.93	1.0	2.29	3.13	4.8	7.92	11.68
OHD-B-ucihar	83.5	10048	243.92	<b>0.3</b>	<b>1.96</b>	<b>1.94</b>	1.5	3.14	3.11	3.9	3.0	4.49	4.45	1.2	1.54	1.53	0.9	1.11	1.10	14.4	7.42	7.38
OHD-M-ucihar	91.9	10048	5.20 <sup>†</sup>	<b>0.6</b>	<b>2.19</b>	<b>2.21<sup>†</sup></b>	2.4	3.42	3.39 <sup>†</sup>	1.9	9.2	7.20	7.20 <sup>†</sup>	0.5	1.09	1.10 <sup>†</sup>	1.2	1.81	1.84 <sup>†</sup>	19.2	14.03	14.05 <sup>†</sup>
LeH-B-ucihar	88.1	5056	123.11	0.0	1.13	1.12	<b>1.0</b>	<b>1.48</b>	<b>1.46</b>	6.7	4.7	1.89	1.88	1.1	1.25	1.25	1.5	1.18	1.17	31.4	9.40	9.42
LeH-M-ucihar	95.3	5056	2.21 <sup>†</sup>	0.0	2.49	2.52 <sup>†</sup>	0.9	3.99	4.20 <sup>†</sup>	8.6	<b>2.1</b>	<b>6.04</b>	<b>4.60<sup>†</sup></b>	5.0	1.72	1.74 <sup>†</sup>	6.1	2.69	2.73 <sup>†</sup>	21.4	36.64	47.13 <sup>†</sup>
LDC-B-ucihar	93.5	256	12.15	<b>0.0</b>	<b>1.62</b>	<b>1.14</b>	0.6	2.31	1.27	4.2	1.8	3.01	1.36	0.0	1.00	0.94	30.4	16.00	1.53	30.4	16.00	1.69
LDC-M-ucihar	94.6	256	70.19	0.0	1.53	1.47	0.0	2.13	2.05	1.8	<b>0.4</b>	<b>2.69</b>	<b>2.59</b>	0.7	1.38	1.32	0.7	1.35	1.31	38.0	16.00	15.46
OHD-B-isolet	85.2	10048	268.56	<b>0.1</b>	<b>1.64</b>	<b>1.59</b>	0.8	2.22	2.14	2.9	2.3	2.95	2.85	0.5	1.37	1.35	0.6	1.02	1.01	23.3	10.57	10.50
OHD-M-isolet	94.4	10048	6.68 <sup>†</sup>	0.0	4.13	4.33 <sup>†</sup>	0.2	6.23	6.16 <sup>†</sup>	1.2	<b>0.9</b>	<b>8.81</b>	<b>8.58<sup>†</sup></b>	1.3	1.29	1.33 <sup>†</sup>	1.0	2.67	2.74 <sup>†</sup>	17.2	42.94	50.58 <sup>†</sup>
LeH-B-isolet	88.5	5056	135.40	0.0	1.10	1.06	0.0	1.28	1.24	1.3	<b>0.0</b>	<b>1.46</b>	<b>1.41</b>	0.9	1.20	1.19	0.7	1.10	1.08	49.3	9.88	9.82
LeH-M-isolet	96.3	5056	2.76 <sup>†</sup>	0.1	2.71	2.76 <sup>†</sup>	<b>0.7</b>	<b>3.95</b>	<b>4.22<sup>†</sup></b>	4.8	1.7	5.37	5.74 <sup>†</sup>	1.0	1.24	1.26 <sup>†</sup>	1.6	1.84	1.87 <sup>†</sup>	36.1	39.50	48.44 <sup>†</sup>
LDC-B-isolet	91.3	256	14.15	0.0	1.27	0.89	<b>0.1</b>	<b>1.67</b>	<b>1.06</b>	4.1	1.0	2.08	1.16	0.0	1.00	0.96	51.8	16.00	1.79	51.8	16.00	1.79
LDC-M-isolet	96.2	256	78.77	0.0	1.41	1.21	<b>0.0</b>	<b>1.92</b>	<b>1.64</b>	3.9	1.6	2.39	2.04	0.8	1.21	1.08	1.0	1.03	0.92	60.6	16.00	15.57
OHD-B-mnist	85.5	10048	301.17	0.1	2.23	2.20	<b>0.3</b>	<b>3.16</b>	<b>3.12</b>	2.3	1.5	4.23	4.18	1.1	2.11	2.09	1.0	1.23	1.22	16.4	8.01	7.95
OHD-M-mnist	96.9	10048	6.25 <sup>†</sup>	0.0	2.28	2.33 <sup>†</sup>	<b>0.4</b>	<b>3.37</b>	<b>3.45<sup>†</sup></b>	2.5	1.9	4.46	4.62 <sup>†</sup>	1.0	1.07	1.07 <sup>†</sup>	1.0	1.24	1.23 <sup>†</sup>	51.2	36.41	39.31 <sup>†</sup>
LeH-B-mnist	93.0	5056	151.09	<b>0.0</b>	<b>1.10</b>	<b>1.08</b>	0.5	1.34	1.32	2.2	2.5	1.62	1.60	0.3	1.05	1.04	0.5	1.04	1.03	40.6	9.82	9.81
LeH-M-mnist	97.7	5056	2.45 <sup>†</sup>	0.0	1.85	1.83 <sup>†</sup>	<b>0.7</b>	<b>2.67</b>	<b>2.63<sup>†</sup></b>	6.0	2.6	3.44	3.42 <sup>†</sup>	1.8	1.09	1.08 <sup>†</sup>	1.8	1.13	1.12 <sup>†</sup>	57.9	39.50	41.56 <sup>†</sup>
LDC-B-mnist	96.2	256	16.51	<b>0.0</b>	<b>1.12</b>	<b>0.97</b>	0.2	1.38	0.99	4.4	1.7	1.67	1.06	0.0	1.00	0.97	61.9	16.00	1.72	61.9	16.00	1.72
LDC-M-mnist	97.5	256	66.41	0.1	1.40	1.31	<b>0.4</b>	<b>1.91</b>	<b>1.78</b>	3.3	1.3	2.37	2.20	1.1	1.10	1.03	1.2	1.04	0.97	56.8	16.00	15.59

**Table 6.** Benchmark names are formatted as {training algorithm}-{HDC type}-{dataset}. For HDC types, B is short for BSC and M is short for MAP. Runtime measured on local laptop is marked as <sup>†</sup>. Other runtimes are collected on a MCU. For each benchmark, the **bold** OMEN execution achieves comparable or better *AL*, better *DRR*, and *SUR* metrics than heuristic termination executions that achieve less than 10% accuracy loss.

Baseline	Encode	Distance Compute	Statistical Tests	Total
UNOPTIMIZED	243.40	0.10	-	243.50
OMEN $\alpha=5\%$	77.67	0.04	0.59	78.30

**Table 7.** Average latency (ms) breakdown for UNOPTIMIZED and OMEN  $\alpha=5\%$  in OHD-BSC-ucihar benchmark.

is because this termination strategy must be paired with the specific training algorithm described in previous work [53].

#### 7.4 Robustness to Hardware Noise

HDC is noise-resilient due to its *fully distributed* property and thus is well-suited for noisy emerging hardware technologies. Provided the noise distribution is identically distributed as hardware errors can be treated as samples from a single probability distribution, OMEN’s assumptions still hold and can be readily applied without adaption or tuning.

We evaluate OMEN with a noisy hardware model in which hypervector bits are flipped with a fixed probability, referred to as the bit error rate (BER). The model approximates the error distribution on the resistive RAM (ReRAM) hardware substrate [29, 49, 69]. ReRAM offers benefits such as larger memory density, nonvolatility, and faster access, but is prone to bit corruption. We investigate BERs of 0.0215 and 0.1273, BERs collected from fabricated ReRAM hardware when setting the memory density to 2 bits per cell (2BPC), and 3BPC, respectively [69]. We evaluate the error resiliency of the baselines on LDC-BSC benchmarks as they have the smallest model and

are suitable for deployment in resource-constrained and potentially error-prone scenarios. We use the same parameters ( $\alpha$  and thresholds) as in the error-free experiments.

The results are summarized in Table 8. OMEN delivers a decent accuracy-performance trade-off with the accuracy loss robustly bounded by  $\alpha$ . Across the 6 benchmarks, OMEN  $\alpha=1\%$  delivers 1.01–1.56 $\times$  dimension reduction with 0.0–0.2% accuracy drop, OMEN  $\alpha=5\%$  delivers 1.07–2.19 $\times$  dimension reduction with 0.3–1.4% accuracy drop, and OMEN  $\alpha=10\%$  delivers 1.16–2.91 $\times$  dimension reduction with 1.2–3.9% accuracy drop.

In contrast, DIFF suffers an accuracy loss of 0.6–5.9%, and MEAN suffers from >20% accuracy loss on all benchmarks. ABSOLUTE becomes over conservative, as the overall similarities between hypervectors decrease under bit corruption, and thus it becomes harder to meet the absolute threshold. ABSOLUTE delivers only 1.02–1.22 $\times$  dimension reduction, while still suffering a 0.6–2.9% accuracy drop. OMEN with  $\alpha=5\%$  consistently outperforms all heuristics that have <10% accuracy loss in both accuracy and dimension reduction. Although tuning the heuristics’ thresholds with bit corruptions could improve accuracy, it is time-consuming, and the error characteristics in deployment may differ from the training environment.

## 8 Related Work

**Global HDC optimizations.** Various theoretical results explored the relationship between the hypervector size and computation accuracy [10, 16, 17, 37, 40, 40, 42, 45, 56, 66, 73].



Benchmark	Unoptimized		Omen $\alpha=1\%$		Omen $\alpha=5\%$		Omen $\alpha=10\%$		Diff		Absolute		Mean	
	acc	# dim	AL	DRR	AL	DRR	AL	DRR	AL	DRR	AL	DRR	AL	DRR
ucihar-2BPC	92.7	256	0.2	1.56	<b>0.8</b>	<b>2.19</b>	2.1	2.91	3.8	1.75	1.5	1.22	28.2	13.84
ucihar-3BPC	88.6	256	0.1	1.15	<b>1.4</b>	<b>1.45</b>	3.9	1.81	5.9	1.31	2.9	1.11	42.9	8.54
isolet-2BPC	91.2	256	0.1	1.21	<b>0.5</b>	<b>1.55</b>	1.9	1.89	2.8	1.20	1.3	1.04	58.4	15.99
isolet-3BPC	80.0	256	0.0	1.02	<b>0.3</b>	<b>1.10</b>	1.2	1.22	3.3	1.08	0.6	1.01	62.4	15.92
mnist-2BPC	94.5	256	0.0	1.09	<b>0.4</b>	<b>1.29</b>	1.5	1.51	0.6	1.04	1.3	1.03	60.7	15.85
mnist-3BPC	75.0	256	0.1	1.01	<b>0.7</b>	<b>1.07</b>	2.3	1.16	1.3	1.03	1.0	1.02	52.6	15.87

**Table 8.** Error resiliency results for BSC-LDC trained benchmarks. Results and **bold** text are formatted the same as in Table 6.

Exploiting these theories, Yi and Achour [72] proposed a static analysis framework that computes optimal hypervector sizes and parameters given the computation specifications. The theories are difficult to apply to practical machine learning, such as classification tasks, and researchers have developed various heuristics for global HDC optimizations in these tasks. Morris et al. traded off hypervector size and accuracy [52] by compressing the class hypervectors. Safa et al. and Imani et al. improved inference performance by sparsifying the class hypervectors [33, 61]. OMEN can be applied on top of most global HDC optimization techniques, as long as they do not break the fully distributedness of hypervectors.

**Per-inference HDC optimizations.** Chuang et al. [9] identified and exploited the difference in the difficulty of inference between test samples to improve HDC performance. They defined the confidence metric as the Hamming distance difference between the classes with the smallest distances and used the binary or integer HDC model depending on whether the confidence surpasses a threshold. Imani et al. [30] used a small 2000-bit model and a large 10000-bit model. They run inference first on the small model and checked whether the similarity value of any class surpasses a specified threshold. If not, the inference used the large model. Similar similarity threshold-based approaches have been used in genome sequence search and adaptive HD model quantization [8, 53, 74]. The confidence metrics in these techniques are heuristics and provide no guarantee on the accuracy loss. OMEN defines confidence as statistical significance and provides accuracy guarantees.

**HDC/VSA.** Hyperdimensional Computation (HDC) or Vector Symbolic Architectures (VSA) is a brain-inspired computational paradigm that has received increasing attention [5, 19, 24, 41, 43, 44]. HDC/VSA has shown great potential in various tasks, such as Genome Sequence Matching [8], Internet of Things systems [31], illness detection [20, 61], Out-of-Distribution Detection [70], image translation [64], neuro-symbolic AI [27], time-series classification [62], and robust hashing [25]. Various training frameworks for HDC/VSA based ML have been proposed, including iterative, online, optimization-based methods [6, 14, 15, 26, 54]. Due to its error resiliency, HDC/VSA is recognized as a suitable computing paradigm for emerging hardware technologies [32, 38, 48, 65]. Besides BSC and MAP that we target in this paper, other

HDC/VSA systems can work with unit-length, real-valued, complex-valued, sparse binary, residue numbers and matrix-as-element hypervectors [1, 18, 47, 55–59, 63, 71]. OMEN’s general idea can be applied to other HDC/VSA systems with slight modifications on the inferential statistics. We leave the application OMEN to other HDC/VSA systems, to other applications, and to combine with other techniques to future work.

**Applying statistical approaches to computing systems.** Unlike heuristic approaches, statistical approaches are theoretically solid and offer guarantees. They are thus widely applied in various computer science problems, including machine learning [34], model checking [51], software testing [35], and data stream processing [2]. To our knowledge, no previous work has focused on edge setting and OMEN is the first statistical approach-based tool to optimize the performance of edge computing systems.

## 9 Conclusion

We presented OMEN, a dynamic HDC optimizer that improves inference performance by terminating the inference early on a per-input basis while providing guarantees on the accuracy loss. OMEN identified good termination points per input. We demonstrated that OMEN achieved a significant inference speedup with only a very small drop in accuracy, and OMEN’s termination improved over heuristic approaches. OMEN could be used to perform an early termination even in the presence of hardware errors. OMEN required little-to-no modification of the HDC training algorithms and could be applied on top of other optimization methods. OMEN’s new statistical perspective on HDC opens doors for applying statistical research to HDC classification problems. OMEN is highly configurable and extendable. Based on OMEN, future work could study optimal strategies for selecting termination points and more aggressive and perhaps unsound HDC performance optimizations.

**Acknowledgements.** We thank Arun Kumar Kuchibhotla for the helpful discussion on exchangeability, Christopher Matthew De Sa for serving as our shepherd and for helpful feedback to the paper, and the anonymous reviewer who took the time to write a review of ~5000 words with detailed comments and actionable advice. This research was supported by the Stanford SystemX Alliance.

## References

- [1] Diederik Aerts, Marek Czachor, and Bart De Moor. Geometric analogue of holographic reduced representation. *Journal of Mathematical Psychology*, 53(5):389–398, 2009.
- [2] Henrique CM Andrade, Buğra Gedik, and Deepak S Turaga. *Fundamentals of stream processing: application design, systems, and analytics*. Cambridge University Press, 2014.
- [3] Davide Anguita, Alessandro Ghio, Luca Oneto, Xavier Parra, Jorge Luis Reyes-Ortiz, et al. A public domain dataset for human activity recognition using smartphones. In *Esann*, volume 3, page 3, 2013.
- [4] Arthur Asuncion and David Newman. UCI machine learning repository, 2007.
- [5] Sercan Aygun, Mehran Shoushtari Moghadam, M Hassan Najafi, and Mohsen Imani. Learning from hypervectors: A survey on hypervector encoding. *arXiv preprint arXiv:2308.00685*, 2023.
- [6] Toygun Basaklar, Yigit Tuncel, Shruti Yadav Narayana, Suat Gumussoy, and Umit Y Ogras. Hypervector design for efficient hyperdimensional computing on edge devices. *arXiv preprint arXiv:2103.06709*, 2021.
- [7] Julius R Blum, H Chernoff, M Rosenblatt, and H Teicher. Central limit theorems for interchangeable processes. *Canadian Journal of Mathematics*, 10:222–229, 1958.
- [8] Hanning Chen, Yeseong Kim, Elaheh Sadredini, Saransh Gupta, Hugo Latapie, and Mohsen Imani. Sparsity controllable hyperdimensional computing for genome sequence matching acceleration. In *2023 IFIP/IEEE 31st International Conference on Very Large Scale Integration (VLSI-SoC)*, pages 1–6. IEEE, 2023.
- [9] Yu-Chuan Chuang, Cheng-Yang Chang, and An-Yeu Andy Wu. Dynamic hyperdimensional computing for improving accuracy-energy efficiency trade-offs. In *2020 IEEE Workshop on Signal Processing Systems (SiPS)*, pages 1–5. IEEE, 2020.
- [10] Kenneth L. Clarkson, Shashanka Ubaru, and Elizabeth Yang. Capacity analysis of vector symbolic architectures, 2023.
- [11] AM Dean and JS Verducci. Linear transformations that preserve majorization, schur concavity, and exchangeability. *Linear algebra and its applications*, 127:121–138, 1990.
- [12] Don Kurian Dennis, Durmus Alp Emre Acar, Vikram Mandikal, Vinu Sankar Sadasivan, Harsha Vardhan Simhadri, Venkatesh Saligrama, and Prateek Jain. Shallow rnns: A method for accurate time-series classification on tiny devices. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, pages 12916–12926, 2019.
- [13] Persi Diaconis and David Freedman. Finite exchangeable sequences. *The Annals of Probability*, pages 745–764, 1980.
- [14] Shijin Duan, Yejia Liu, Shaolei Ren, and Xiaolin Xu. Lehd: Learning-based hyperdimensional computing classifier. In *Proceedings of the 59th ACM/IEEE Design Automation Conference*, pages 1111–1116, 2022.
- [15] Shijin Duan, Xiaolin Xu, and Shaolei Ren. A brain-inspired low-dimensional computing classifier for inference on tiny devices. *arXiv preprint arXiv:2203.04894*, 2022.
- [16] E Paxon Frady, Denis Kleyko, and Friedrich T Sommer. A theory of sequence indexing and working memory in recurrent neural networks. *Neural Computation*, 30(6):1449–1513, 2018.
- [17] Stephen I Gallant and T Wendy Okaywe. Representing objects, relations, and sequences. *Neural computation*, 25(8):2038–2078, 2013.
- [18] Ross W Gayler. Multiplicative binding, representation operators & analogy (workshop poster). 1998.
- [19] Lulu Ge and Keshab K Parhi. Classification using hyperdimensional computing: A review. *IEEE Circuits and Systems Magazine*, 20(2):30–47, 2020.
- [20] Lulu Ge and Keshab K Parhi. Applicability of hyperdimensional computing to seizure detection. *IEEE Open Journal of Circuits and Systems*, 3:59–71, 2022.
- [21] Dirk Goldhahn, Thomas Eckart, Uwe Quasthoff, et al. Building large monolingual dictionaries at the leipzig corpora collection: From 100 to 200 languages. In *LREC*, volume 29, pages 31–43, 2012.
- [22] Hui Han and Julien Siebert. TinyML: A systematic review and synthesis of existing research. In *2022 International Conference on Artificial Intelligence in Information and Communication (ICAIC)*, pages 269–274. IEEE, 2022.
- [23] Yizeng Han, Gao Huang, Shiji Song, Le Yang, Honghui Wang, and Yulin Wang. Dynamic neural networks: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(11):7436–7456, 2021.
- [24] Eman Hassan, Yasmin Halawani, Baker Mohammad, and Hani Saleh. Hyper-dimensional computing challenges and opportunities for ai applications. *IEEE Access*, 10:97651–97664, 2021.
- [25] Mike Heddes, Igor Nunes, Tony Givargis, Alexandru Nicolau, and Alex Veidenbaum. Hyperdimensional hashing: A robust and efficient dynamic hash table. In *Proceedings of the 59th ACM/IEEE Design Automation Conference*, pages 907–912, 2022.
- [26] Alejandro Hernández-Cano, Namiko Matsumoto, Eric Ping, and Mohsen Imani. OnlineHD: Robust, efficient, and single-pass online learning using hyperdimensional system. In *2021 Design, Automation & Test in Europe Conference & Exhibition (DATE)*, pages 56–61, 2021.
- [27] Michael Hersche, Mustafa Zeqiri, Luca Benini, Abu Sebastian, and Abbas Rahimi. A neuro-vector-symbolic architecture for solving raven’s progressive matrices. *Nature Machine Intelligence*, 5(4):363–375, 2023.
- [28] Sture Holm. A simple sequentially rejective multiple test procedure. *Scandinavian journal of statistics*, pages 65–70, 1979.
- [29] ER Hsieh, M Giordano, B Hodson, A Levy, SK Osekowsky, RM Radway, YC Shih, W Wan, TF Wu, X Zheng, et al. High-density multiple bits-per-cell 1t1r RRAM array with gradual set/reset and its effectiveness for deep learning. In *2019 IEEE International Electron Devices Meeting (IEDM)*, pages 35–6. IEEE, 2019.
- [30] Mohsen Imani, Chenyu Huang, Deqian Kong, and Tajana Rosing. Hierarchical hyperdimensional computing for energy efficient classification. In *Proceedings of the 55th Annual Design Automation Conference*, pages 1–6, 2018.
- [31] Mohsen Imani, Yeseong Kim, Behnam Khaleghi, Justin Morris, Haleh Alimohamadi, Farhad Imani, and Hugo Latapie. Hierarchical, distributed and brain-inspired learning for internet of things systems. In *2023 IEEE 43rd International Conference on Distributed Computing Systems (ICDCS)*, pages 511–522, 2023.
- [32] Mohsen Imani, Abbas Rahimi, Deqian Kong, Tajana Rosing, and Jan M Rabaey. Exploring hyperdimensional associative memory. In *2017 IEEE International Symposium on High Performance Computer Architecture (HPCA)*, pages 445–456. IEEE, 2017.
- [33] Mohsen Imani, Sahand Salamat, Behnam Khaleghi, Mohammad Samragh, Farinaz Koushanfar, and Tajana Rosing. SparseHD: Algorithm-hardware co-optimization for efficient high-dimensional computing. In *2019 IEEE 27th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM)*, pages 190–198. IEEE, 2019.
- [34] Gareth James, Daniela Witten, Trevor Hastie, Robert Tibshirani, et al. *An introduction to statistical learning*, volume 112. Springer, 2013.
- [35] Ramesh Johari, Pete Koomeen, Leonid Pekelis, and David Walsh. Peeking at a/b tests: Why it matters, and what to do about it. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1517–1525, 2017.
- [36] Rakhee Kallimani, Krishna Pai, Prasoon Raghuwanshi, Sridhar Iyer, and Onel LA López. TinyML: Tools, applications, challenges, and future research directions. *Multimedia Tools and Applications*, 83(10):29015–29045, 2024.
- [37] Pentti Kanerva. Fully distributed representation. *PAT*, 1(5):10000, 1997.
- [38] Geethan Karunaratne, Manuel Le Gallo, Giovanni Cherubini, Luca Benini, Abbas Rahimi, and Abu Sebastian. In-memory hyperdimensional computing. *Nature Electronics*, 3(6):327–337, 2020.
- [39] Michael Klass and Henry Teicher. The central limit theorem for exchangeable random variables without moments. *The Annals of Probability*, pages 138–153, 1987.

- [40] Denis Kleyko, Connor Bybee, Ping-Chen Huang, Christopher J Kymn, Bruno A Olshausen, E Paxon Frady, and Friedrich T Sommer. Efficient decoding of compositional structure in holistic representations. *Neural Computation*, 35(7):1159–1186, 2023.
- [41] Denis Kleyko, Mike Davies, Edward Paxon Frady, Pentti Kanerva, Spencer J Kent, Bruno A Olshausen, Evgeny Osipov, Jan M Rabaey, Dmitri A Rachkovskij, Abbas Rahimi, et al. Vector symbolic architectures as a computing framework for emerging hardware. *Proceedings of the IEEE*, 110(10):1538–1571, 2022.
- [42] Denis Kleyko, Evgeny Osipov, Alexander Senior, Asad I Khan, and Yaşar Ahmet Şekercioglu. Holographic graph neuron: A bioinspired architecture for pattern processing. *IEEE transactions on neural networks and learning systems*, 28(6):1250–1262, 2016.
- [43] Denis Kleyko, Dmitri Rachkovskij, Evgeny Osipov, and Abbas Rahimi. A survey on hyperdimensional computing aka vector symbolic architectures, part ii: Applications, cognitive models, and challenges. *ACM Computing Surveys*, 55(9):1–52, 2023.
- [44] Denis Kleyko, Dmitri A Rachkovskij, Evgeny Osipov, and Abbas Rahimi. A survey on hyperdimensional computing aka vector symbolic architectures, part i: Models and data transformations. *ACM Computing Surveys (CSUR)*, 2021.
- [45] Denis Kleyko, Antonello Rosato, Edward Paxon Frady, Massimo Panella, and Friedrich T. Sommer. Perceptron theory can predict the accuracy of neural networks. *IEEE Transactions on Neural Networks and Learning Systems*, pages 1–15, 2023.
- [46] Arun Kumar Kuchibhotla. Exchangeability, conformal prediction, and rank tests. *arXiv preprint arXiv:2005.06095*, 2020.
- [47] Christopher J Kymn, Denis Kleyko, E Paxon Frady, Connor Bybee, Pentti Kanerva, Friedrich T Sommer, and Bruno A Olshausen. Computing with residue numbers in high-dimensional representation. *ArXiv*, 2023.
- [48] Jovin Langenegger, Geethan Karunaratne, Michael Hersche, Luca Benini, Abu Sebastian, and Abbas Rahimi. In-memory factorization of holographic perceptual representations. *Nature Nanotechnology*, 18(5):479–485, 2023.
- [49] Binh Q Le, Akash Levy, Tony F Wu, Robert M Radway, E Ray Hsieh, Xin Zheng, Mark Nelson, Priyanka Raina, H-S Philip Wong, Simon Wong, et al. RADAR: A fast and energy-efficient programming technique for multiple bits-per-cell RRAM arrays. *IEEE Transactions on Electron Devices*, 68(9):4397–4403, 2021.
- [50] Yann LeCun, Corinna Cortes, and CJ Burges. Mnist handwritten digit database. *ATT Labs [Online]*. Available: <http://yann.lecun.com/exdb/mnist>, 2, 2010.
- [51] Axel Legay, Benoît Delahaye, and Saddek Bensalem. Statistical model checking: An overview. In *International conference on runtime verification*, pages 122–135. Springer, 2010.
- [52] Justin Morris, Mohsen Imani, Samuel Bosch, Anthony Thomas, Helen Shu, and Tajana Rosing. CompHD: Efficient hyperdimensional computing using model compression. In *2019 IEEE/ACM International Symposium on Low Power Electronics and Design (ISLPED)*, pages 1–6. IEEE, 2019.
- [53] Justin Morris, Si Thu Kaung Set, Gadi Rosen, Mohsen Imani, Baris Aksanli, and Tajana Rosing. AdaptBit-HD: Adaptive model bitwidth for hyperdimensional computing. In *2021 IEEE 39th International Conference on Computer Design (ICCD)*, pages 93–100. IEEE, 2021.
- [54] Nuntipat Narkthong, Shijin Duan, Shaolei Ren, and Xiaolin Xu. MicroVSA: An ultra-lightweight vector symbolic architecture-based classifier library for always-on inference on tiny microcontrollers. In *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2*, pages 730–745, 2024.
- [55] T.A. Plate. Holographic reduced representations. *IEEE Transactions on Neural Networks*, 6(3):623–641, 1995.
- [56] Tony A Plate. *Distributed representations and nested compositional structure*. Citeseer, 1994.
- [57] Tony A Plate. Analogy retrieval and processing with distributed vector representations. *Expert systems*, 17(1):29–40, 2000.
- [58] Tony A Plate. Holographic reduced representation: Distributed representation for cognitive structures. 2003.
- [59] D.A. Rachkovskij. Representation and processing of structures with binary sparse distributed codes. *IEEE Transactions on Knowledge and Data Engineering*, 13(2):261–276, 2001.
- [60] Dmitriy A Rachkovskiy, Sergey V Slipchenko, Ernst M Kussul, and Tatyana N Baidyk. Sparse binary distributed encoding of scalars. *Journal of Automation and Information Sciences*, 37(6), 2005.
- [61] Ali Safa, Ilja Ocket, Francky Catthoor, and Georges Gielen. Support-HDC: Hyperdimensional computing with scalable hypervector sparsity. In *Proceedings of the 2023 Annual Neuro-Inspired Computational Elements Conference, NICE '23*, page 20–25, New York, NY, USA, 2023. Association for Computing Machinery.
- [62] Kenny Schlegel, Peer Neubert, and Peter Protzel. HDC-MiniROCKET: Explicit time encoding in time series classification with hyperdimensional computing. In *2022 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2022.
- [63] Javier Snaider and Stan Franklin. Modular composite representation. *Cognitive Computation*, 6:510–527, 2014.
- [64] Justin Theiss, Jay Leverett, Daeil Kim, and Aayush Prakash. Unpaired image translation via vector symbolic architectures. In Shai Avidan, Gabriel Brostow, Moustapha Cissé, Giovanni Maria Farinella, and Tal Hassner, editors, *Computer Vision – ECCV 2022*, pages 17–32, Cham, 2022. Springer Nature Switzerland.
- [65] Simon Thomann, Hong L. G. Nguyen, and Hussam Amrouch. HW/SW codesign for approximate in-memory computing. In *2022 23rd International Symposium on Quality Electronic Design (ISQED)*, pages 1–6, 2022.
- [66] Anthony Thomas, Sanjoy Dasgupta, and Tajana Rosing. Theoretical foundations of hyperdimensional computing. *Journal of Artificial Intelligence Research*, 72:215–249, 2021.
- [67] Yaman Umuroglu, Nicholas J Fraser, Giulio Gambardella, Michaela Blott, Philip Leong, Magnus Jahre, and Kees Vissers. FINN: A framework for fast, scalable binarized neural network inference. In *Proceedings of the 2017 ACM/SIGDA international symposium on field-programmable gate arrays*, pages 65–74, 2017.
- [68] Abraham Wald. Sequential tests of statistical hypotheses. In *Breakthroughs in statistics: Foundations and basic theory*, pages 256–298. Springer, 1992.
- [69] Anjiang Wei, Akash Levy, Pu (Luke) Yi, Robert Radway, Priyanka Raina, Subhasish Mitra, and Sara Achour. PBA: Percentile-based level allocation for multiple-bits-per-cell RRAM. In *ICCAD*, 2023.
- [70] Samuel Wilson, Tobias Fischer, Niko Sünderhauf, and Feras Dayoub. Hyperdimensional feature fusion for out-of-distribution detection. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 2644–2654, January 2023.
- [71] Calvin Yeung, Zhuowen Zou, and Mohsen Imani. Generalized holographic reduced representations. *arXiv preprint arXiv:2405.09689*, 2024.
- [72] Pu (Luke) Yi and Sara Achour. Hardware-aware static optimization of hyperdimensional computations. *Proc. ACM Program. Lang.*, 7(OOPSLA2), oct 2023.
- [73] Tao Yu, Yichi Zhang, Zhiru Zhang, and Christopher M De Sa. Understanding hyperdimensional computing for parallel single-pass learning. *Advances in Neural Information Processing Systems*, 35:1157–1169, 2022.
- [74] Zhuowen Zou, Hanning Chen, Prathyush Poduval, Yeseong Kim, Mahdi Imani, Elaheh Sadredini, Rosario Cammarota, and Mohsen Imani. BioHD: an efficient genome sequence search platform using hyperdimensional memorization. In *Proceedings of the 49th Annual International Symposium on Computer Architecture*, pages 656–669, 2022.