
HYPERDIMENSIONAL VECTOR TSETLIN MACHINES WITH APPLICATIONS TO SEQUENCE LEARNING AND GENERATION

Christian D. Blakely

Centre for Artificial Intelligence Research
University of Agder
Grimstad, Norway
christian.blakely@uia.no

August 2024

ABSTRACT

We construct a two-layered model for learning and generating sequential data that is both computationally fast and competitive with vanilla Tsetlin machines, adding numerous advantages. Through the use of hyperdimensional vector computing (HVC) algebras and Tsetlin machine clause structures, we demonstrate that the combination of both inherits the generality of data encoding and decoding of HVC with the fast interpretable nature of Tsetlin machines to yield a powerful machine learning model. We apply the approach in two areas, namely in forecasting, generating new sequences, and classification. For the latter, we derive results for the entire UCR Time Series Archive and compare with the standard benchmarks to see how well the method competes in time series classification.

1 Introduction

A large part of any design of a data learning agent is in feature extraction of the underlying data, and how it is computed and represented. The best processes for extracting features for learning information from data typically take advantage of expert knowledge of the underlying data to either expose the most relevant features, reduced noise, and extract the most amount of independent information in the data. For many types of datasets, this might be challenging due to factors such as incoherence, abstractedness, or the sheer amount of noise present in the data. In designing features for Tsetlin machines, one is tasked to booleanize (or binarize) the underlying data, and under the presence of noise, this can be challenging. Furthermore, for notoriously complex high-dimensional data like noisy sequences, graphs, images, signal spectra, and natural language, creating encodings that are also interpretable for human reasoning in any post-hoc process can be difficult due to creating logic AND expressions that both take advantage of the relevant information in the data, but also lead to accurate expressions that can compete with other machine learning models.

In this paper, we explore using Hyperdimensional Vector Computing (HV computing, or simply HVC) as an input layer to a novel Tsetlin machine architecture and apply it to learning, classifying, predicting, and generating sequences. Here, we argue that HVC can provide a robust layer of feature extraction due to the many computational advantages. This approach was first introduced in [1] and here, we streamline the approach to focus on sequences while further leveraging other attributes of HCV such as N-Gram sequence encoding and associative memory, while combining with TMs, to create a powerful hybrid methodology while remaining minimalist in memory sizes of the overall model.

1.1 Hyperdimensional Vector Representations

There are many hyperdimensional vector (HV) representations that have been introduced for various applications over the past few decades. Combining with Tsetlin machines, the most natural approach is to use what is called in the literature as Binary Spatter Codes (BSC), or simply binary vectors of high dimension. These high-dimensional vectors, often referred to as HVs, are typically composed of thousands of bits. HVC is inspired by the way the human brain

processes information, providing a robust and efficient means for computing and reasoning. Given the fact that Tsetlin machines learn directly from logical expressions on binarized data, BSC provides a natural input layer for Tsetlin machines.

While HVs are not typically treated as elements of a traditional vector space over a field, they share some properties with vector spaces. The algebra of binary HVs, which we will give an overview of in section 2 is characterized by:

- **Bundling:** Combining multiple vectors into one using majority voting or similar methods.
- **Binding:** Associating vectors using the XOR operation.
- **Unbinding:** Retrieving original vectors from bound vectors using the XOR operation.
- **Similarity Measurement:** Using Hamming distance to measure the similarity between vectors.
- **Perturbation:** Slightly modifying vectors to represent related but distinct concepts.

These operations form the core algebraic framework for HVC, enabling robust and efficient computation in binary spaces

1.2 Motivation and contributions

One of the biggest attractions of HVC is the ability to represent virtually any type of data through the use of the vector operations. The main motivation in this paper is to build a strategy for learning and generating sequences by establishing a robust encoding from multidimensional sequences to HVs. We also desire that the encoding is then naturally suited for learning with TMs, and so we will provide empirical evidence of an attractive coupling of both HVC and TMs.

This paper contributes a new approach to sequence learning, where we show we can compete with current SOTA results on fundamental benchmarks. Furthermore, we also show how combining some of the features of HVC with TMs can produce a highly versatile forecasting model for many types of sequences, while generating new sequences of any length with are similar in Hamming distance to any given sequence. Lastly, we will show that the models have a very light footprint in memory, making it an attractive approach for online learning in small embedded systems.

1.3 Organization of Paper

We first give a brief overview of HVC along with the core algebraic framework we will be using throughout the paper. We then discuss the encoding strategy we will be using for various types of sequences, along with designing a so-called associative memory which will be very useful in sequence generation and forecasting.

The next section briefly outlines the TM architecture we will be using throughout, with references to more in-depth treatments of the approach.

We will present an in-depth numerical treatment of the approach by first outlining the hybrid structure of our HVC-TM architecture, and then providing an algorithm for time series learning and forecasting. In the final section of the paper, we will provide results for applying our proposed approach on the entire UCR Time Series archive for classifying many types of time series.

Finally, we will conclude the study with some alternatives in design and further next steps with possible applications.

2 Review of Binary HVs

We will leverage heavily throughout our approach some fundamentals of HVs. A fantastic survey and introduction to HVC can be found in the recent monograph [2], which also gives insights into current research trends. To give a quick overview on why one would want to use HVs

- **Robustness:** HVs are resilient to noise and errors. Small perturbations in the vectors do not significantly affect their overall similarity, making HVC robust to noisy data.
- **Scalability:** HVC can easily handle high-dimensional data, and the operations on HVs are computationally efficient.
- **Simplicity:** The operations on HVs, such as bundling, binding, and similarity measurement, are straightforward and can be implemented efficiently.
- **Flexibility:** HVs can represent a wide range of data types and structures, from simple scalars to complex symbolic representations.

The operations we use of binary HVs differ from standard real vectors in that we don't use multiplication and addition. Instead, we apply operations defined by binding and bundling on HVs, followed by the inclusion and notion of a distance for comparing HVs. Perturbations of vectors are also used to represent the notion of ordering in HVs (for example, vector **A** comes before **B** in a sequence). Here we give a brief introduction to these operations

2.1 Bundling

Bundling is the operation of combining multiple HVs into a single HV. This operation is analogous to the concept of superposition in physics, where multiple states are combined. In HVC, bundling is typically implemented using element-wise majority voting and will be represented throughout the paper as a $+$.

Given binary HVs **A**, **B**, and **C**, the bundled HV **S** can be expressed as $\mathbf{S} = \text{Majority}(\mathbf{A}, \mathbf{B}, \mathbf{C})$ where for each bit position i :

$$S_i = \begin{cases} 1 & \text{if } A_i + B_i + C_i \geq 2 \\ 0 & \text{otherwise.} \end{cases} \quad (1)$$

Because of the fact we take a majority, it is important that we only apply this operation on an odd number of HVs.

2.2 Binding

Binding is the operation of combining two HVs to form a new HV that represents their association. In HVC, binding is typically implemented using the element-wise XOR operation. Given binary HVs **A** and **B**, the bound HV **C** is: $\mathbf{C} = \mathbf{A} \oplus \mathbf{B}$ where \oplus denotes the XOR operation for each bit position i , $C_i = A_i \oplus B_i$. The unbinding operation is used to retrieve an original HVs that was bound together with another, and is frequently used in decoding HVC representations. This is achieved by applying the XOR operation with one of the original HVs to the bound HV.

Given the bound HV **C** and one of the original HVs (say, **A**), we can retrieve the other HV (**B**) as follows: $\mathbf{B} = \mathbf{C} \oplus \mathbf{A}$. This works because the XOR operation is its own inverse. Specifically, $\mathbf{A} \oplus \mathbf{A} = \mathbf{0}$ (the zero vector) giving $\mathbf{C} \oplus \mathbf{A} = (\mathbf{A} \oplus \mathbf{B}) \oplus \mathbf{A} = \mathbf{A} \oplus \mathbf{A} \oplus \mathbf{B} = \mathbf{0} \oplus \mathbf{B} = \mathbf{B}$.

Thus, by XORing the bound HV **C** with one of the original HVs (**A**), we effectively cancel out **A** and retrieve the other original HV (**B**). For example, suppose we have two binary HVs **A** and **B**: $\mathbf{A} = [1, 0, 1, 1]$, $\mathbf{B} = [0, 1, 0, 1]$. Their bound HV **C** is $\mathbf{C} = \mathbf{A} \oplus \mathbf{B} = [1 \oplus 0, 0 \oplus 1, 1 \oplus 0, 1 \oplus 1] = [1, 1, 1, 0]$. To retrieve **B**, we unbind **C** with **A**: $\mathbf{B} = \mathbf{C} \oplus \mathbf{A} = [1 \oplus 1, 1 \oplus 0, 1 \oplus 1, 0 \oplus 1] = [0, 1, 0, 1]$. We will often compare HVs, and this is done via a similarity measure called the Hamming distance which is applied between two binary vectors **A** and **B** of equal length is defined as the number of positions at which the corresponding bits differ, namely $d_H(\mathbf{A}, \mathbf{B}) = \sum_{i=1}^n (A_i \oplus B_i)$. The Hamming distance is a measure of dissimilarity between two binary vectors, so a Hamming distance of 0 indicates that the vectors are identical, while a larger Hamming distance indicates greater dissimilarity.

The final operation useful for sequences is perturbation operator. It is the operation of slightly modifying a HV to represent a related but distinct vector. We implement this by simply shifting cyclically the bits of the vector j places. We will denote a perturbation applied to a vector hv by $p(hv, j)$. The inverse operation of perturbation ip is then just $p(hv, -j)$. This will be used for representing the j th position in a sequence. We will typically call the vector resulting in applying a j perturbation on a vector as a position vector.

2.3 Sequence operations

We can now apply these operations for representing and processing sequences of data. Several recent studies have been proposed for encoding sequences, such as in [3] and [4]. In a similar fashion, the goal of our approach will be to encode a time series sequence $\mathbf{s} \in \mathbb{R}^n$ into a binary vector dimension of D . In this paper, we will represent sequences (or time series) as a stationary sequence of scalars $\mathbf{s} = [s_1, s_2, \dots, s_n]$. We will show how we represent these as HVs, leveraging the operations of bundling, binding, and perturbation to capture the structure and relationships within the sequence.

The first step is to introduce what we call an Interval embedding. Given a quantization of dimension Q , since we assume our sequences are stationary, they are bounded and thus have a minimum and maximum interval, $[m_0, m_1]$. We divide the interval $[m_0, m_1]$ into Q buckets, can represent each bucket by an HV (we will denote these as hv). Thus the process of mapping $\mathbf{S} = [s_1, s_2, \dots, s_n]$ into a collection of HVs $E_S = [hv_1, hv_2, \dots, hv_n]$ will be represented as E . This is shown in the following algorithm.

Algorithm 1: Interval Embedding

Input: Scalar value x , Interval range $[m_0, m_1]$, Number of divisions Q
Output: Hyperdimensional vector hv

```

1 Initialize  $step \leftarrow \frac{m_1 - m_0}{Q}$ ;
2 Initialize  $intervals \leftarrow \{\}$ ;
3 for  $i \leftarrow 0$  to  $Q - 1$  do
4   |  $intervals[i] \leftarrow$  Random hyperdimensional vector;
5 end
6  $index \leftarrow \left\lfloor \frac{x - m_0}{step} \right\rfloor$ ;
7  $hv \leftarrow intervals[index]$ ;
8 return  $hv$ ;
```

And now we use the interval embedding function to encode the original sequence into a sequence of HV interval embeddings.

Algorithm 2: Encoding Scalar Sequence

Input: Scalar sequence $S = [s_1, s_2, \dots, s_n]$, Interval embedding function E
Output: Encoded sequence $E_S = [hv_1, hv_2, \dots, hv_n]$

```

1  $E_S \leftarrow \{\}$ ;
2 for  $i \leftarrow 1$  to  $n$  do
3   |  $hv_i \leftarrow E(s_i)$ ;
4   |  $E_S \leftarrow E_S \cup \{hv_i\}$ ;
5 end
6 return  $E_S$ ;
```

With any sequence now being transformed into a sequence of HVs, the next step is encoding this sequence into one HV. Here we apply both perturbations and binding with N-Grams. To encode the original sequence $E_S = \{hv_1, hv_2, hv_3, \dots, hv_n\}$, apply a convolution comprised of the perturbation representing a position vector applied to the hv_i followed by a binding with a Gram vector $p(\cdot, j) \oplus G_j$. If the number of Gram vectors is, for example 3, the first element would be $GV_2 = (p(hv_0, 0) \oplus G_0) \oplus (p(hv_1, 1) \oplus G_1) \oplus (p(hv_2, 2) \oplus G_2)$. This can be summarized below.

Binding with Position Vectors

- Generate random position vectors for each position in the N-Gram $G_j, j = 1, \dots, N$.
- Bind each element of the N-Gram with its corresponding position vector using the XOR operation.

For short, we will define the N-Gram vector at the t -th temporal index by $GV_t = (p(hv_t, 0) \oplus G_0) \oplus (p(hv_{t-1}, 1) \oplus G_1) \oplus (p(hv_{t-N}, N) \oplus G_N)$ With the position encoded, and the N-gram contains the spatial information of N subsequent samples with different timestamps, making it a spatiotemporal HV.

- Divide the sequence into overlapping N-Grams. For example, for a sequence $E_S = \{hv_1, hv_2, hv_3, \dots, hv_n\}$ and $N = 3$, the N-Grams are $\{hv_1, hv_2, hv_3\}, \{hv_2, hv_3, hv_4\}$, etc.
- Apply the permutation operator $p(\cdot, j)$ to the j th position, yielding a position vector $P_j := p(hv, j)$
- At each position, bind the position vector with the j th Gram HV
- Bind all the N-Gram position vectors together

Once we have generated GV_0, \dots, G_{L-N-1} , where L is the length of the sequence and N is the number of grams, then we can then apply a bundling operation to combined the elements in the HV space. Recall that bundling the bound HVs of the N-Grams using element-wise majority voting. In summery, the steps of encoding a sequence is shown if the figure below.

The first step in the top encodes all the scalar values of the entire sequence of time series values into a sequence of HVs based on the interval embedding procedure. In order to encode these HVs into a single HV, we apply N-Gram encoding by first applying a permutation to the respective position giving encoding position, and then we proceed by apply XOR operation with the N-Gram vectors. This is essentially a convolution across the time series sequence given not with temporal and spatial information. With the collection $L - N - 1$ HVs post XORing with N-Gram vectors, we then

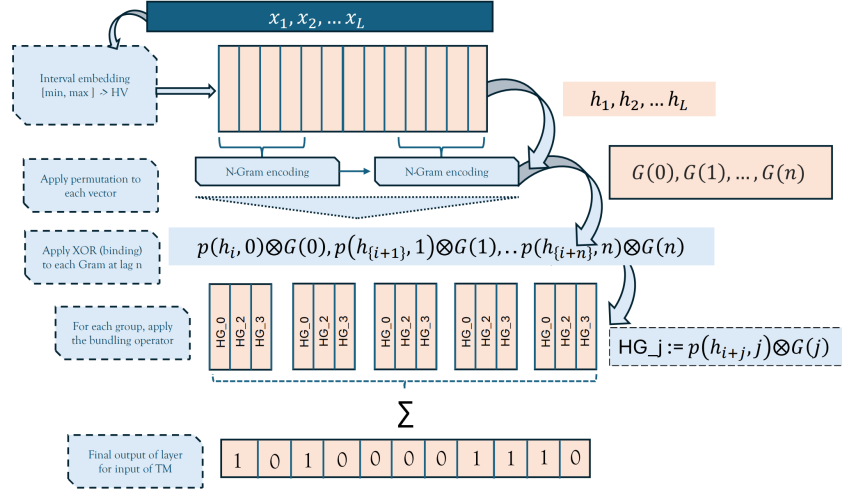


Figure 1: The procedure for encoding a sequence into a hypervector.

apply a bundle operation to store all the information in one HV. This vector is now ready as input into a TM model, along with a class label or the next value in the sequence (forecasting).

We now apply this procedure to build an associative memory for a collection of sequences that will be used for our hybrid forecasting approach.

2.4 Associative Memory

Associative memory in the context of HVC refers to a memory model where stored items are retrieved based on their similarity to a query item. It leverages the properties of high-dimensional spaces to perform efficient and robust retrieval operations. We propose a hybrid model of forecasting that is a multistep step approach based on associative memory.

To build the associative memory for a given set of sequences, the process involves taking all possible N-Grams of the series and encoding them into HVs, while storing them with their associative values in the sequence. More formally, for a set of M sequences $\mathbf{S} = \{S_0, \dots, S_M\}$, where $S_i := s_0, s_1, \dots, s_L$, we extract all possible N-Grams from each sequence. We will define

$$\mathbf{H}_t^m := (p(hv_t^m, 0) \oplus G_0) \oplus (p(hv_{t-1}^m, 1) \oplus G_1) \oplus (p(hv_{t-N}^m, N) \oplus G_N) \quad (2)$$

by the t th N-Gram HV from the m th sequence, and then define the mapping

$$\mathcal{M}(\mathbf{H}_t^m) \mapsto s_t^m,$$

for all sequence values s_t^m . This creates a dictionary, namely a (key, value) pairing where \mathcal{M} is the Memory, the key is a HV, and the value is the interval embedding HV representing a scalar. The process of querying this mapping is called retrieving memory from the associated memory. To retrieve an item similar to a query HV \mathbf{Q} , the memory HV \mathbf{M} is queried, and the stored HV \mathbf{H} that is most similar to \mathbf{Q} is retrieved. The similarity is measured using the Hamming distance:

$$\mathbf{H}_{\text{retrieved}} = \arg \min_{\mathbf{H} \in \text{Memory}} d_H(\mathbf{Q}, \mathbf{H}) \quad (3)$$

Given a set of stored HVs $\{\mathbf{H}_1, \mathbf{H}_2, \dots, \mathbf{H}_k\}$, the memory HV \mathbf{M} is: $\mathbf{M} = \sum_{i=1}^k \mathbf{H}_i$. To retrieve the item most similar to the query \mathbf{Q} , we calculate the Hamming distance between \mathbf{Q} and any proposed \mathbf{H}_i . This process allows for efficient and robust retrieval of stored items based on their similarity to the query item. Before we describe the forecasting and classification approach of HVTMs, we first review the TM architecture we will be using.

3 Tsetlin Machine Architecture

In this paper we employ a recently developed architecture for TM learning originally proposed in [5]. The architecture shares an entire pool of clauses between all output classes, while introducing a system of weights for each class type. The learning of weights is based on increasing the weight of clauses that receive a Type Ia feedback (due to true

positive output) and decreasing the weight of clauses that receive a Type II feedback (due to false positive output). This architectural design allows to determine which clauses are inaccurate and thus must team up to obtain high accuracy as a team (low weight clauses), and which clauses are sufficiently accurate to operate more independently (high weight clauses). The weight updating procedure is given in more detail in [5]. Here we illustrate the overall scheme in Figure 2. Notice that each clause in the shared pool is related to each output by using a weight dependant on the output class and the clause. The weights that are learned during the TM learning steps and are multiplied by the output of the clause for a given input. Thus clause outputs are related to a set of weights for each output class.

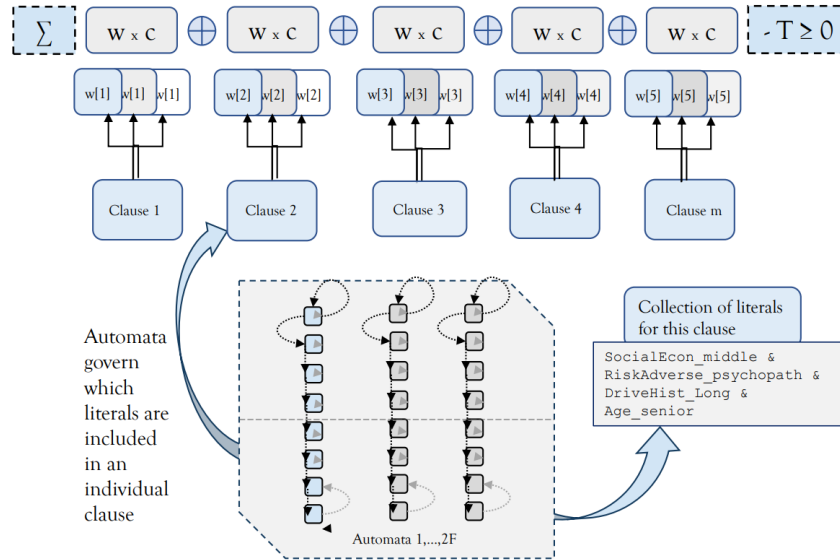


Figure 2: Showing a collection of clauses and their relationship with weights that are learned during TM training. Each clause contains a set of literals that are also learned during feedback.

An additional parameter we use for our TM architecture is to ensure that a maximum number of literals is allowed for every clause. This has the effect of ensuring that the clauses do not become too "dense", proposed in [6].

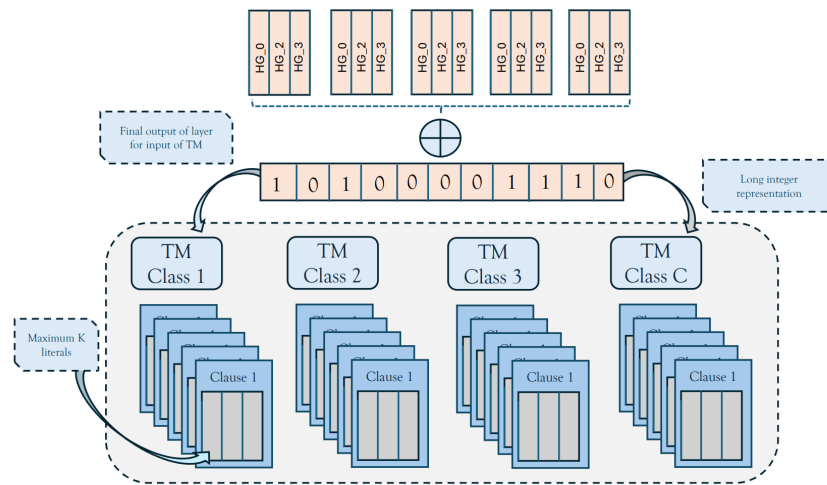


Figure 3: Input into the TM is the N-Gram encoded latest N values of the sequence. This input is used to predict the next value governed by the Q -class TM

4 Time Series Classification

In this section we do an in-depth study on how the HVTM performs in classifying many types of time series. For this study, we apply our proposed methodology on the UCR Time Series Classification archive [7] which is a widely used comprehensive collection of datasets for bench-marking time series classification algorithms. Boasting 128 different time series datasets from a broad range of domains, including medicine, biology, finance, motion tracking, and sensor data with a diverse range of time series lengths, some as short as a few dozen with very few training samples or having a high imbalance of a certain class, all while having strong benchmark tracking, it is the premier data archive for testing new time series classification models.

Here we apply our methodology to all 128 data sets, and for this we fix the HV strategy, namely setting

- HV dimension = 5000
- N-Gram length = 3
- Quantization = 50

For the TM architecture, we randomly choose 20 different configurations where

- number clauses between [100, 2000]
- number max literals per clause [10, 100]
- specificity between [10f, 20f]
- threshold between [0, 100].

For each random configuration, we run on 10 Epochs, and take the accuracy after the final Epoch. Finally, to measure the final performance, we then take the top 10 model configurations and report the maximum accuracy.

4.1 Numerical Results

In order to compare our approach with the published benchmarks of [7], for all 128 time series, we compare the accuracy of our classification with the best accuracy recorded in [7]. Their approach compares Euclidean distance, and two Dynamic Time Warping (DTW) computations (one with a learned/optimized parameter warping parameter, and one with a fixed parameter).

To visually compare, we plot the pair (accuracy_{HVTM} , accuracy_{UCR}).

Without having chance for optimization on controlling the dimension of the HV, nor the number of Grams used to encode the hypervector, and using a random parameterization for the TM, we can see that the approach has performed quite well simply using the "out-of-the-box" default settings. The HVTM method improves or competes in accuracy (cutoff at 2 percent) of the optimal benchmark provide by [7] in roughly 78 percent of the data sets.

To give an overview of the performance comparison, we employ a scatter plot showing the accuracy of the HVTM approach on each data set versus the DTW benchmark. In Figure 4, the scatter plot shows the clustering along $x = y$ line, demonstrating that the approach is very competitive with the benchmark.

To get even more insights into the (out)performance, we break down further the results into four categories: data type, length of time series, number of training samples available, and number of classes in order to gain possible information on where HVTMs outperformance different training setups. Figure 5 shows the comparison of performance across the different data types and the length of the time series. The first two bars in each category represent the mean accuracy for both HVTM and the DTW benchmark. The third bar show the percentage of times HVTM outperformed the benchmark in terms of accuracy. We highlight (using a shadow on the bars) the categories in which HVTM outperformed DTW at least 60 percent of the time.

We can see that HVTM outperformed the DTW benchmark on most of the data set types, including Motion, Images, and ECGs. A few categories, including EOG, and Traffic, only two data sets fall into these category making performance comparisons difficult.

Next we see that if we compare the performances broken down into time series length, there is clear dominance by HVTM over most lengths. However, it seemed to struggle with very short series, namely series of length 24-80. Furthermore, for mid-range series, 277-500 in length, there is no clear outperformance. However, for series greater than 500 in length there seems to be a clear dominance by the HVTM approach.

We also are interested in how the methods compare given the number of training samples. More training samples can yield more dispersion of information for each class, but here HVTM has no problem in outperforming DTW

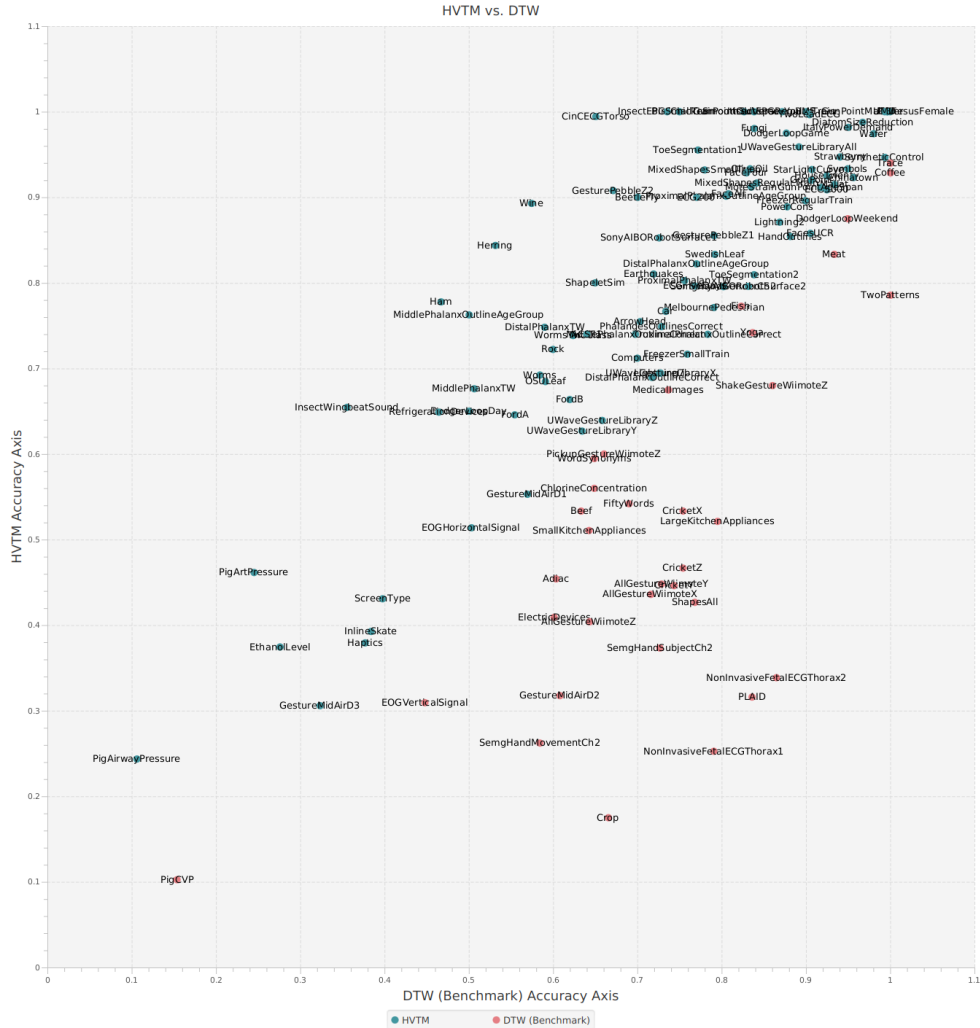


Figure 4: Scatter plot showing the accuracy of the performance on each data set versus the benchmark. Here the X-axis is the benchmark accuracy and the Y-axis is the HVTM accuracy

benchmarks for a very few number of training samples. However, DTW seems to outperform for a very large number of training samples (more than 1000). In regards to number of classes, HVTM can handle less than 10 classes very well, and outperforms the DTW approach fairly well. However, for many more classes, the HVTM approach with the given HV parameters seems to struggle. This could be simply due to only utilizing a dimension of 5000 for the HV system.

In most categories across all 4 bar charts of performance comparisons, there does not exist a large discrepancy between average error rates. The only instances of any large discrepancy is for a large number of classes. To improve on this, some further optimization on either the dimension D of the HV space, the number of N-Grams, or allowing more candidate TM models might improve the accuracy of the HVTM approach. One could also consider encoding more information into the HV, such as the class label (or potential class label) and first training purely in the HV space by using associative memory.

4.2 Hybrid Predictions

With the computational framework for both HVC and our TM architecture, we now show a straightforward approach to forecasting and generating new sequences based on encoded historical sequences. Prediction involves using the encoded HV of a sequence to predict the next element(s) in the sequence. Here we make a few assumptions on the underlying time series data

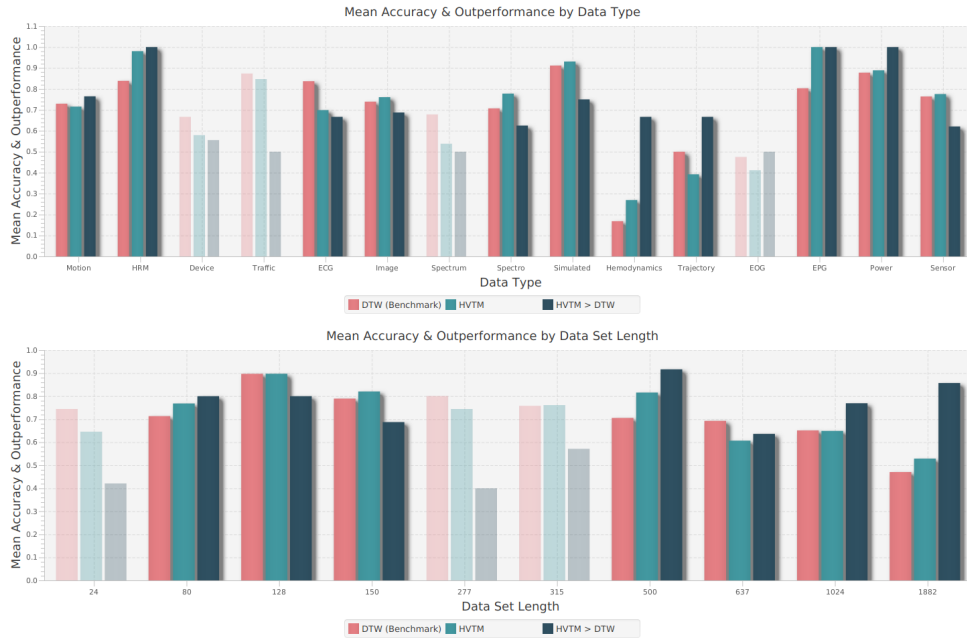


Figure 5: Mean accuracy and (out)performance results into four categories: data type, length of time series



Figure 6: Mean accuracy and (out)performance results into four categories: number training samples, number classes

- The time series data is stationary, namely we can assume there is a lower and upper bound indefinitely during any time range we consider
- While only one time series is necessary and sufficient, any additional time series for training purposes can be considered from the the same data generating process
- No NAs/missing values exist in the time series, namely all the data for training is available

We also note that we make no assumption on the uniformity of the data sampling of the underlying time series. So all series are considered to have uniform time stamps. However, some applications could utilize the distance between each time stamp as a feature, but we will not consider this in the study.

Our forecasting approach considers two steps for training: 1) encoding the time series into a dictionary of N-Grams mapped to time series values by creating an associative memory, and 2) Training a TM on the dictionary of N-Gram HVs labeled with the final sequence value in that N-Gram.

- Given the latest observation s_n in the time series, find the next observation \hat{s}_{n+1} that has the highest similarity from the associative memory based on the current context s_{n-N}, \dots, s_n
- Use associative memory to retrieve the most similar past context and predict the next element based on this context.
- For the proposed observation \hat{s}_{n+1} , recursively generate then next $\hat{s}_{n+2}, \hat{s}_{n+3}, \dots, \hat{s}_{n+K}$
- For each proposed observation \hat{s}_{n+1} , using the TM with Q classes, where Q is the quantization level, predict the each next observation \bar{s}_{n+1+i} given the gram vector G
- Combine the forecasts $\hat{s}_{n+1} + \bar{s}_{n+1}$

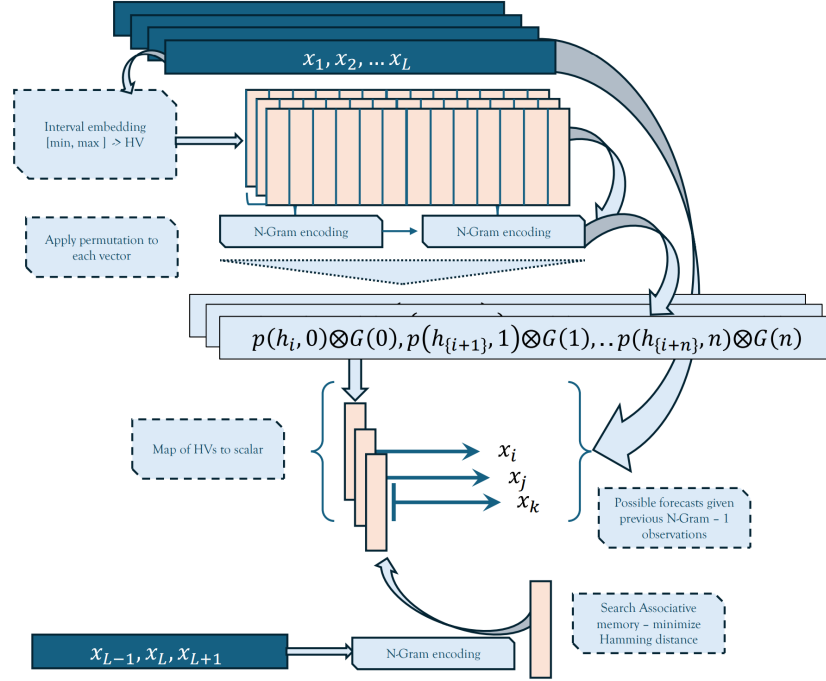


Figure 7: Showing the steps for encoding based on N-Gram encoding, then forecasting given the current encoded vector based on the two most recent observations.

The second step is done by using as input into the HVTM the N-Gram encoded vector of the last N values of the time series, followed by predicting the next quantized value in the sequence. This gives us a second forecast, from which we can derive a weighted average from the associative memory forecast.

This is shown in Figure 8. Each block of window size M leads to a group of HVs that is the collection of N-Gram encoding of length N . This HV is used as the data with the label being the quantized value of the next value in the sequence. This procedure produces many labeled samples from which we can train the TM using the HVs. For prediction, we combine the prediction made by using the latest HV in the time series sequence and combine it with the prediction made by the TM model.

4.2.1 Numerical Experiments

In order to empirically determine how well the forecasting and sequence generation performs, we begin by simulating deterministic sequences and comparing different length forecasts with the simulated sequence. For this, our model

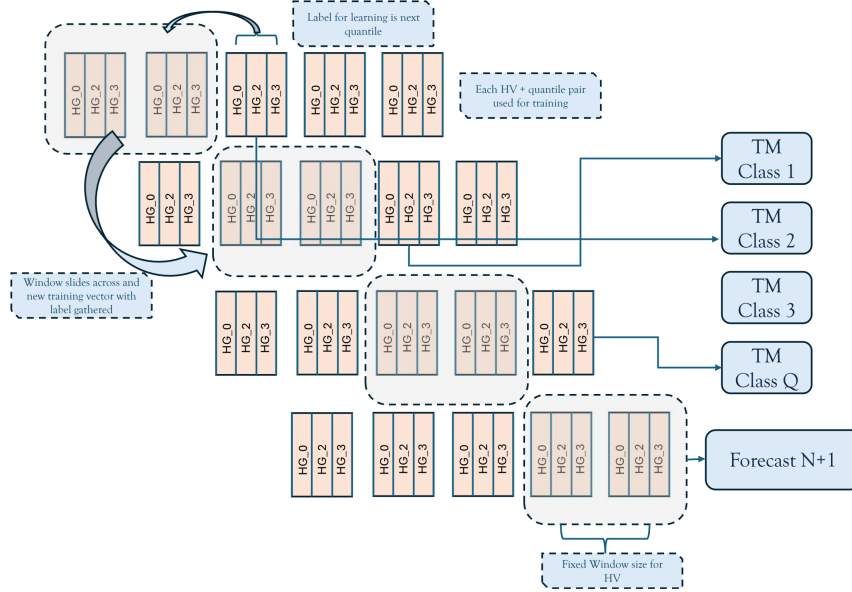


Figure 8: Input into the TM is the N-Gram encoded latest N values of the sequence. This input is used to predict the next value governed by the Q -class TM

is a harmonic series $s(t) = a \sin(t2\pi) * \cos(bt2\pi + ct) + \sin(2\pi + t)$ where we randomly choose $a, b, c \in [0, 1]$ to generate a sequence. For our model setup, we again use $D = 5000$ and the number of N-Grams we set at $N = 5$. For the TM setup we chose the number of clauses to be 1000, the threshold at 100, and the specificity to 15. The max number of literals per clause was fixed at 50. These choices were made given the performances from the classification numerical experiments. They generally yielded top results.

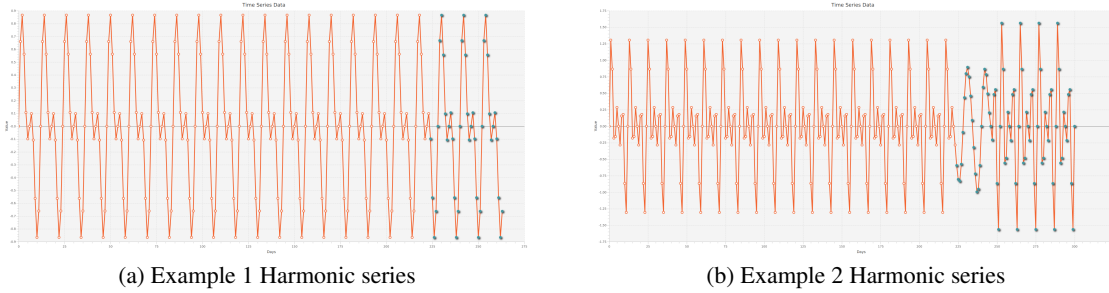


Figure 9: Generating 24 steps ahead for two different deterministic series.

Figures 9 and 10 show 4 examples of various harmonic series where we sample 220 observations. For each series the first 220 observations were encoded using N-Gram windowing technique shown in figure 7 after which we then learned the samples from the encoded N-Grams and labeled them with the next values in the sequence, as shown in figure 8.

We can see here that the hybrid forecasts are able to extract the harmonic structure well enough to generalize a reproduction of the periodicity involved.

Each new value generated is only dependent on the previous 4 values, since the N-Gram length is 5. Subsequent forecasts are then used to generate new values in the sequence.

Next we expand on this by now simulating series from stochastic models. Here we consider a simple AR(1) model with very high AR coefficient, an ARMA(1,1) model. Finally in order to see how well it performs on seasonal data, we also sample data from a seasonal AR model with strong seasonal coefficients.

A seasonal model produces new challenges not seen in AR(MA) models as the sequence generator also has to "learn" what frequency the seasonality occurs at. We can see in 13 that is able to find estimate the seasonality fairly well.

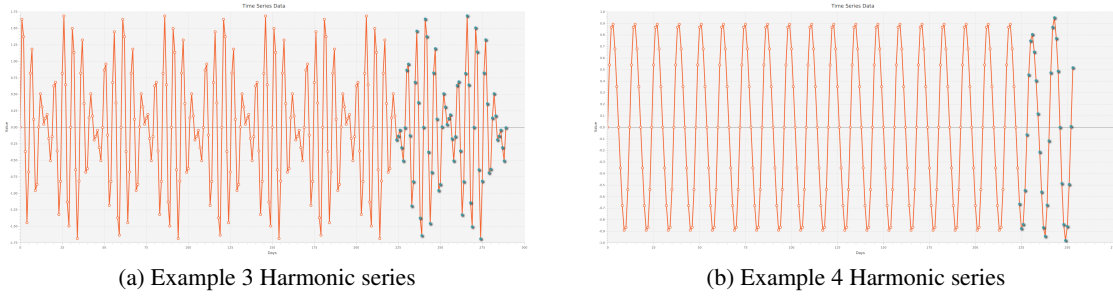


Figure 10: Generating 24 steps ahead for two different deterministic series.

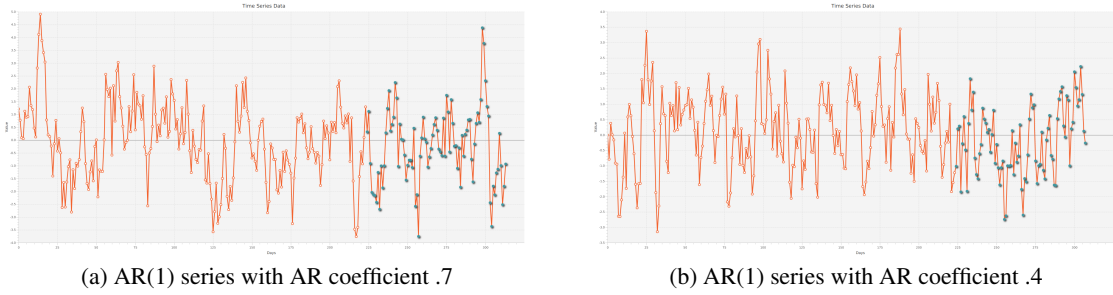


Figure 11: Generating 24 steps or more ahead for two different stochastic series.

To generate empirical results on the forecasting ability of the approach, we sample from each model 50 times and generate sequentially a 24-step ahead forecast. For each forecast, we compute the error with the ground truth next value in that sequence. The time series have been normalized between $[-1, 1]$ to ensure normalized regular record keeping of the errors (so for example, in an extreme case, a mean error of $\sqrt{(48)/24}$ would imply all ground truth values were 1 but forecasts were -1, which gives an upper bound on the errors.)

Model	Coefficients	N-Gram 3	N-Gram 5	N-Gram 7
Harmonics		0.078 (0.015)	0.044 (0.012)	0.051 (0.022)
AR(1)	0.2	0.298 (0.020)	0.102 (0.020)	0.143 (0.052)
AR(1)	0.4	0.245 (0.114)	0.152 (0.085)	0.215 (0.131)
AR(1)	0.7	0.245 (0.011)	0.140 (0.081)	0.200 (0.121)
SAR(1)	0.1; 0.7	0.511 (0.315)	0.312 (0.181)	0.415 (0.317)
SAR(1)	0.3; 0.7	0.421 (0.273)	0.552 (0.121)	0.651 (0.323)
ARMA(1,1)	0.7; -0.1	0.245 (0.014)	0.152 (0.074)	0.201 (0.027)

Table 1: Table of Model Coefficients and N-Gram Values with corresponding standard errors

We can see that the forecasting performs best on the deterministic Harmonic series, which is to be expected. the standard error is near 7 percent over a 24 step forecast, for 3 N-Gram encoding, and around 4 percent for 5 N-Gram encoding. In general, the best performance comes from 5 N-Gram encoding which is able to capture best the periodicity. For the AR models, we see they are fairly consistent regardless of the coefficient, with a rather tight standard deviation.

The seasonal AR model poses more challenges because not only does it have to correctly forecast the next value, but also determine what what frequency the seasonality is occurring. Naturally, with the larger variance at seasonal periods, the standard error will typically be larger. For both seasonal models, it still performs quite well with a "best" performance at around 31 percent error. Lastly, for the ARMA model the forecasting performance seems to be nearly as good as the standard AR model with a lower standard error near 15 percent.

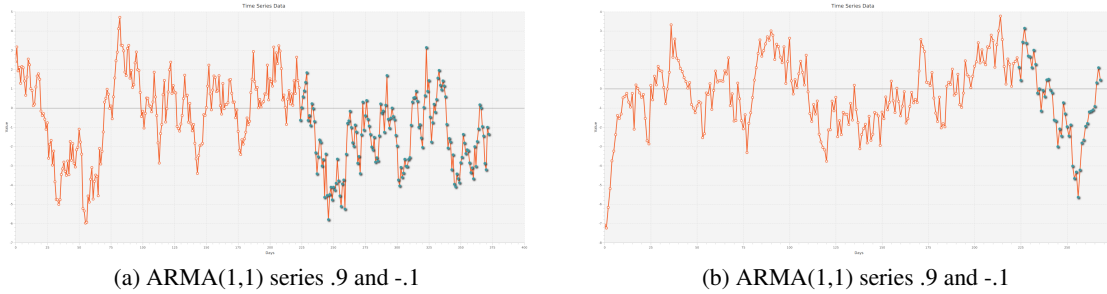


Figure 12: Generating 24 steps or more ahead for two different stochastic series.

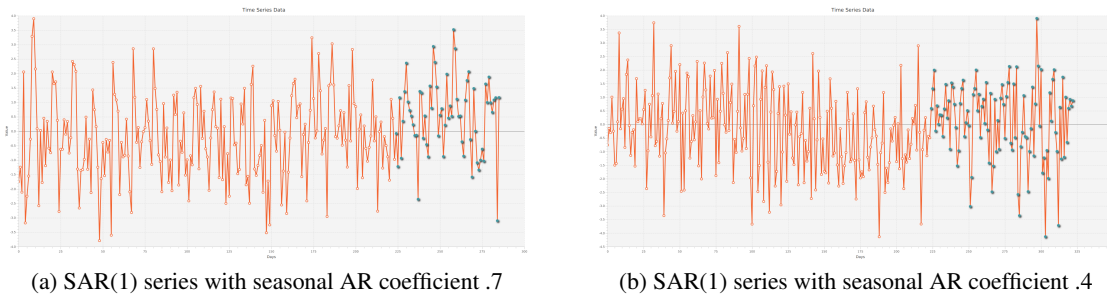


Figure 13: Generating 24 steps or more ahead for two different stochastic series.

5 Conclusion

In this paper we proposed a straightforward approach to combining the strengths of HVC with the machine learning power of Tsetlin machines to give a flexible approach for learning sequences, classifying them, and generating new ones. Our approach relied heavily on encoding procedures for sequences where we took advantage of a powerful N-Gram structure for spatialtemporal learning. Since the resulting HVs can be decoded as well, combined with the innate interpretable structure of Tsetlin machines, this hybrid approach could be an attractive alternative to larger models such as Deep learning sequence models (LSTMs, RNNs) which rely heavily on large scale weight estimation via matrix vector multiplication, backpropagation, optimization and parameter hypertuning.

5.1 Model Size

This leads us to an additional note that we have not discussed in the paper but would be valuable for future work. Model size and computation learning in terms of energy expenditure could be optimized and improved greatly by taking advantage of the fact we are using binary vectors. The size of a model based on HVs depends on several factors, including the dimensionality of the vectors, the number of vectors used in the model (number of N-Grams and interval embedding vectors), but we can safely conclude that reasonably, with a dimension of 10k for or HV system, and with a thousand vectors, we have a size of 1,220KBs total ($10,000\text{bits}/8=1,250\text{bytes}$ times 1000). Expanding to even larger dimensions, say $D = 100\text{k}$, we still only need a few MBs of chip memory. Now coupling with the TM, say 1000 clauses, each with max 32 bits for each automata, we are looking at a trivial amount of additional memory.

5.2 Future work

Many direction exist in terms of where these models could be challenged next. The first area of improvement would be in understanding better the underperformance of some of the classification results, including when many classes are present. Here, it is of the author's hypothesis that simply a higher dimension, along with a more careful pruning and selection of model criterion could be achieved. Perhaps even varying the number of N-Grams, to see if they play a large role for a larger number of classes. Since no optimization was done in choosing the TM parameters, this would also be an area to look at more. Since the goal was to see if the approach could work "straight-out-of-the-box", we wanted to have the least amount of tinkering as possible in model architecture setups.

The second area of work would be into taking a deeper dive into forecasting real-world series with more challenging dynamics. For example, some macroeconomic data, central orderbook data, and asset returns from the financial world, or signal sequences from the medical field.

A third area of work would be to take into account additional time series to model and forecast multivariate time series. Here, the attraction would be to incorporate more features from either the underlying series (for example, volatility, seasonality), and see how classification and forecasting can be tackled when additional features are included.

References

- [1] Vojtech Halenka, Ahmed K. Kadhim, Paul F. A. Clarke, Bimal Bhattarai, Rupsa Saha, Ole-Christoffer Granmo, Lei Jiao, and Per-Arne Andersen. Exploring effects of hyperdimensional vectors for tsetlin machines, 2024.
- [2] 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*, 55(6):1–40, December 2022.
- [3] Kenny Schlegel, Peer Neubert, and Peter Protzel. Hdc-minirocket: Explicit time encoding in time series classification with hyperdimensional computing, 2022.
- [4] Dmitri A. Rachkovskij and Denis Kleyko. Recursive binding for similarity-preserving hypervector representations of sequences, 2022.
- [5] Sondre Glimsdal and Ole-Christoffer Granmo. Coalesced multi-output tsetlin machines with clause sharing, 2021.
- [6] Sebastian Østby, Tobias M. Brambo, and Sondre Glimsdal. The sparse tsetlin machine: Sparse representation with active literals, 2024.
- [7] Hoang Anh Dau, Eamonn Keogh, Kaveh Kamgar, Chin-Chia Michael Yeh, Yan Zhu, Shaghayegh Gharghabi, Chotirat Ann Ratanamahatana, Yanping, Bing Hu, Nurjahan Begum, Anthony Bagnall, Abdullah Mueen, Gustavo Batista, and Hexagon-ML. The ucr time series classification archive, October 2018.