# Non-volatile Hierarchical Temporal Memory: Hardware for Spatial Pooling

Lennard G. Streat[*], Dhireesha Kudithipudi[†], Kevin Gomez[§]

Nanocomputing Research Laboratory, Rochester Institute of Technology, Rochester, NY 14623[*†]

Seagate Technology, Shakopee, MN 55379[§]

Email: lgs8331@rit.edu[*], dxkeec@rit.edu[†],kevin.gomez@seagate.com[§]

*Abstract*—**Hierarchical Temporal Memory (HTM) is a biomimetic machine learning algorithm imbibing the structural and algorithmic properties of the neocortex. Two main functional components of HTM that enable spatio-temporal processing are the spatial pooler and temporal memory. In this research, we explore a scalable hardware realization of the spatial pooler closely coupled with the mathematical formulation of spatial pooler. This class of neuromorphic algorithms are advantageous in solving a subset of the future engineering problems by extracting non-intuitive patterns in complex data. The proposed architecture, Non-volatile HTM (NVHTM), leverages large-scale solid state flash memory to realize a optimal memory organization, area and power envelope. A behavioral model of NVHTM is evaluated against the MNIST dataset, yielding 91.98% classification accuracy. A full custom layout is developed to validate the design in a TSMC 180nm process. The area and power profile of the spatial pooler are 30.538mm$^2$ and 64.394mW, respectively. This design is a proof-of-concept that storage processing is a viable platform for large scale HTM network models.**

## I. Introduction

**H**IERARCHICAL temporal memory (HTM) is a biomimetic machine learning algorithm, designed with the aim of capturing key functional properties of the mammalian brain's neocortex to solve pattern recognition problems. HTM theory was originally proposed by Jeff Hawkins in [1], and subsequently formalized in [2], [3]. Since its inception, HTM has evolved as a machine intelligence algorithm with its core being a cortical learning algorithm (CLA) [4], [5], [6]. Given spatio-temporal data, HTM can use learned representations to perform a type of time-dependent regression[7]. Few applications that have shown promise in this area include predicting taxi passenger counts [8], and anomaly detection in streaming data.

In its most recent theoretical formulation, HTM consists of two main functional units: a spatial pooler(SP) and a temporal memory (TM). These two substructures model the spatial and temporal patterns within data, respectively. The SP is responsible for mapping an input data on to a sparse distributed representation and the TM is responsible for learning sequences and making predictions. When combined, they form what is referred to as a "*region*". An HTM region, depicted in Fig. 1, consists of multiple columns (SP). Each column consists of four cells (TM). These cells use proximal dendrites to connect the feedforward input, via synapses, to a column. There are distal segments which connect cells within a region. Regions may be connected hierarchically to form
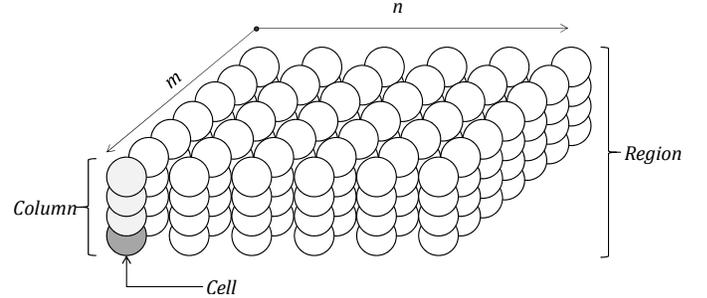


Fig. 1. Example HTM region, consisting of a $6 \times 7$ column space ($n \times m$); each column is composed of 4 cells. Columnar computation is handled by the spatial pooler, and cellular-level operations occur within the temporal memory.

larger systems. A region may be theoretically n-dimensional in size. For more structural details of HTM, please refer to [9].

The focus of this research is on the SP, which maps encoded input (referred to as an "*input region*") into a higher-dimensional space. Output generated by an SP region is an N-dimensional sparse binary vector or a sparse distributed representation (SDR) [10]. Spatial pooling is decomposed into three phases: namely, overlap, inhibition, and learning. Overlap computes the number of active-connected synapses on a column's proximal segment: the set of synapses mapping the input region onto the column. A synapse is considered to be connected if it's permanence (or weight) is above a fixed threshold and active if it is connected to an active input. Inhibition prunes the set of columns with proximal segments above the overlap threshold to obey the SDR constraint. Finally, learning rules are employed to enforce fairness amongst columns (boosting), and to train proximal segments to model patterns in the input region.

CLA was discussed in the literature, with emphasis placed on acceleration methods [11], [12], noise resilience [7], [10], in addition to anomaly detection applications [13]. Despite these efforts, hardware implementations of CLA are still scarce: the same is true for Zeta [14]. This is due in part to volatility of the algorithmic formulation. HTM CLA is studied, in part, in [6], [15]. In this work emphasis is placed on studying the existing CLA models for hardware realization.

Following Numenta's official CLA white paper [6], HTM has seen growth in research interest. For instance, a multi-core software-based CLA was presented in [12]. Zhou and Lou designed a C-based many-core implementation that lever-

aged the Adapteva hardware platform [11]. In the work of Liddiard et al., SP was implemented and evaluated using pseudo-sensory data [16]. However, the analysis was lacking with regard to specificity in the methodology; a non-standard dataset was presented and not characterized in depth. Vyas and Zaveri presented a verilog implementation of SP [17]. More recently, Zyarah designed a 2D mesh architecture that modeled many concepts of CLA and resulted in significant speedup over a matlab-based software implementation [18]. These architectures share some commonalities, and primarily focused on performance scaling, but discussions regarding hardware scalability are lacking. More recently, CLAASIC, a 16-core packet-switched 2D torus network was developed to model CLA [19]. However, as with the previously mentioned models, online learning was not proposed.

In this research, we explore a nonvolatile memory-centric design for realizing the spatial pooler. As dynamic random access memory (DRAM) is reaching its scaling limits, semiconductor technology roadmap projections show that it will likely be replaced by other emerging nonvolatile memories such as spin-transfer torque magnetoresistive random access memory, resistive random access memory, or phase-change memory. Therefore, it is timely to investigate a nonvolatile memory centric processing approach for these systems [20], [21], [22]. Moreover in the future, it is expected that large information processing systems will move from message passing interface to memory-centric programming models, as network latencies improve. In the current market, a non-volatilve memory (NVM), such as flash, provides the optimal balance of cost, performance and energy for the data store in this data-centric architecture. To increase compute parallelism, concurrency (more parallel threads), and minimize data movement and co-locate HTM close to memory, we study flash within solid-state drives (SSD).

In general, this research falls under a broader umbrella of neuromorphic computing. The thesis of modern neuromorphic computing is that if the key principles of brain functionality can be formulated and emulated in a hardware platform, we will be capable of solving a subset of the challenging future engineering problems by extracting non-intuitive patterns in complex data. Exploiting non-volatile memory designs for these systems is a natural choice. Few explorations of neuromorphic algorithms were shown in TrueNorth [23], NeuroGrid [24], and BrainScales [25], architectures, where the synaptic memory is placed adjacent to the neuron circuits. Overall chip architecture will have multiple neuro-synaptic cores which communicate using different asynchronous and synchronous processing. However, the proposed HTM spatial pooler architecture takes a different approach, where the hardware design is custom tuned for the spatial pooler algorithm while exploiting the nonvolatile memory.

The aim of this work is to explore a pragmatic architecture for the HTM spatial pooler and evaluate its ability to scale with respect to classification performance, memory limitations, latency, area, and power consumption. The remainder of this work is structured as follows: In Section II, a hardware-amenable mathematical formulation of the CLA SP is described. Microarchitecture for a storage processor unit (SPU)

implementation, namely Non-volatile HTM (NVHTM), is discussed in Section III. NVHTM is evaluated for its scalability, based on various design constraints in Section IV. Conclusions and future work may be found in Section V.

## II. SPATIAL POOLER MODEL

HTM SP is a three-phase biomimetic unsupervised clustering algorithm that seeks to create a mapping from an encoded input to an SDR. In order, it consists of overlap, inhibition, and learning phases. These operations are computed by individual nodes, or columns, within a two-dimensional column space (Fig. 1). Each aspect of the spatial pooling process is progressively discussed, starting with a simple description of overlap, eventually arriving at the full embodiment of SP with a learning rule. Limitations of the model are discussed to identify where NVHTM diverges from CLA.

### A. Phase I: Overlap

HTM columns are nodes in the two-dimensional column space that contain a proximal segment: synaptic state datastructure. Overlap describes the degree of correspondence between a given proximal segment and the current feedforward stimulus: the input region. This operation is a dot product between the proximal segment and the input region. Columnar proximal segments may be represented as a vector of real-valued permanences, bounded between zero and unity, inclusive. Once the permanence magnitude is above a threshold, $P_{th}$, a potential synapse is promoted to the connected state. Synaptic connectivity is determined by the condition (1); $\vec{c}_i$ represents the proximal segment vector for the $i^{th}$ column and $\vec{C}_i$ is the binary vector of columnar synapse state: a value of 1 is indicative of a connection.

$$\vec{C}_i = \vec{c}_i \geq P_{th} \tag{1}$$

There is a one-to-one mapping between elements within the column vector and the input region. Given that proximal segments may be represented as length-$K$ vectors, the entire column space may be envisioned as a $K \times N$ matrix, $\mathbf{C}$, in which the $i^{th}$ column vector corresponds to the proximal segment of the $i^{th}$ HTM column. Each of the $N$ columns are capable of being connected to a maximum of $K$ synapses.

$$\alpha'_i = \vec{C}_i \cdot \vec{X}_t \tag{2}$$

$$\vec{\alpha'} = \mathbf{C}^T \cdot \vec{X}_t \tag{3}$$

Overlap for the $i^{th}$ column, the dot product of its proximal segment state vector and the binary input vector, is modeled by (2). The entire network may be evaluated in parallel by computing a dot product between the transpose of $\mathbf{C}$ and the current input state vector, $\vec{X}_t$, as shown in (3); $\alpha'_i$ and $\vec{\alpha'}$ are overlap for the $i^{th}$ column and the entire network, respectively. This relationship is extended to account for the boost factor in addition to minimum overlap, yielding

$$\alpha_i = \begin{cases} \alpha'_i \beta_i, & \& \alpha'_i \geq A_{th} \\ 0, & \& otherwise \end{cases} \tag{4}$$

where $\alpha_i$ is the overlap value for the $i^{th}$ column after boosting and minimum overlap thresholding; $A_{th}$ is the minimum overlap threshold; and $\beta_i$ is the boosting factor. In [6], boost factor is restricted to values greater than or equal to unity.

### B. Phase II: Inhibition

A competitive learning process, referred to as inhibition, is leveraged to select candidate columns to enter the active state–thus contributing to the output SDR. Columns are inhibited by neighbors within their inhibition radius and are pruned if their overlap is below that of the $K^{th}$ largest within their radius, which is defined with respect to column-space coordinates (Fig. 1). Inhibition is classified into two forms, namely global and local; the latter of which being significantly more complex to compute.

Inhibition applies a k-select operation and pruning step to remove all columns within a set and an overlap value below that of the selected column. This may be modeled by sorting all columns, $j$, within the inhibition radius of another column, $i$; hereto referred to as $\Lambda_j^i$. Pruning is modeled in (5) by removing columns below a specified index from the aforementioned set, which now becomes the SP output SDR.

$$\mathcal{A}_i = \alpha_i \geq \vartheta[i, d_{th}] \tag{5}$$

$$\vartheta[i, d_{th}] = \Lambda_j^i[min(d_{th}, length(\Lambda_j^i))] \tag{6}$$

where $\mathcal{A}_i$ is the binary inhibited columnar activation; $d_{th}$ is the desired local activity; $\vartheta$ is a bounded indexing function–used to index into $\Lambda_j^i$, subject to $d_{th}$; $\vartheta[i, d_{th}]$ calculates the minimum local activity for the $i^{th}$ column; the $min$ function is used to restrict the range of $\vartheta$ to indexes within $\Lambda_j^i$.

Global inhibition is a special case of the more generalized local model. Given that the inhibition radius of each column is a known constant, computational complexity may be simplified. The global case is configuration in which the "*local*" inhibition radii of all columns are large enough to encapsulate the entire column space: a known fixed quantity. Consequently, the equation governing both global and local cases are, in principle, the same. Implementations of SP in the literature utilize global over local inhibition, yielding comparable results [7]. Furthermore, the global case is more practical for hardware implementation, requiring fewer parameters to be tracked, such as inhibition radii, implicitly removing the requirement for updating these parameters dynamically.

### C. Phase III: Learning

*1) Active Segment Update:* The third, and final, phase of spatial pooling is to train the proximal segment state matrix. SP learning is influenced by network state (due to prior training and observed patterns) in addition to the current feedforward input vector. The current input is dictated by

$$\vec{C}_i^*[j] = \begin{cases} \vec{C}_i[j] + \lambda \vec{X}_t[j] - P_{dec}, & \&\mathcal{A}_i = 1 \\ \vec{C}_i[j], & \&otherwise \end{cases} \tag{7}$$

where $\vec{C}_i^*[j]$ is the new value for the $j^{th}$ permanence within the $i^{th}$ proximal segment state vector; bounded between 0 and 1, inclusive; $\lambda$ is the sum of $P_{inc}$ and $P_{dec}$–the permanence increment and decrement, respectively. Synaptic state is only updated for active columns and is based upon the binary input region vector. If the bit was active, the permanence is incremented, otherwise it is decremented.

*2) Duty Cycle Update:* Following active segment update, duty cycle parameters are updated to enforce fairness between columns. Duty cycle is used as a measure of the extent to which a column participates in representing the encoded input region; it is modeled by an active duty cycle, $D_A$, and overlap duty cycle , $D_O$. $D_A$ is a measure of how frequently a column wins in the inhibition phase; $D_O$ is a measure of the average pre-inhibited overlap for the column in question. Active duty cycle is updated such that

$$D_A^*[i] = \frac{D_A[i](\tau_D - 1) + \mathcal{A}_i}{\tau_D} \tag{8}$$

where $D_A^*[i]$ is the next value for the active duty cycle, $D_A$, of the $i^{th}$ column; $\tau_D$ is the period over which the duty cycle is evaluated and consequently averaged, and $\mathcal{A}_i$ is the current active state for the column. A similar process is repeated for the overlap duty cycle in (9).

$$D_O^*[i] = \frac{D_O[i](\tau_D - 1) + \alpha_i}{\tau_D} \tag{9}$$

$D_O^*[i]$ is the next value for the overlap duty cycle, $D_O$. $\alpha_i$ is the new post-inhibition overlap resulting from the current feedforward input. Although active and overlap duty cycles were implemented having the same period. This is not an inherent limitation of CLA.

*3) Weak Column Boosting:* Weak column boosting, described by (10), seeks to increase the number of synapses that are in the connected state for potentially starved columns. The permanence values of a column, $\vec{C}_i$ are increased by 10% of the permanence threshold [6]. Other such magic numbers were used in the white paper, but further work should seek to optimize the HTM to ascertain optimal parameters using parametric optimization techniques such as simulated annealing.

$$\vec{C}_i^*[j] = \begin{cases} \vec{C}_i[j] + P_{th}/10, & D_O[i] < \widetilde{D}_O[i], \vec{C}_i[j] > 0 \\ \vec{C}_i[j], & otherwise \end{cases} \tag{10}$$

where $\widetilde{D}_O$ is the minimum overlap duty cycle. If the duty cycle of the column in question is below $\widetilde{D}_O$ and the synapse has not already desaturated to a value of 0 (pruned), the permanence values are incremented.

*4) Boost Factor Update:* The overall goal of boosting is to reduce column starvation, promoting an increase in the number of columns contributing to the process of creating a model of the input data. This aspect of the learning algorithm is modeled as

$$\beta_i^* = \begin{cases} \frac{1 - \beta_{max}}{\widetilde{D}_A} D_A[i] + \beta_{max}, & D_A[i] < \widetilde{D}_A \\ \beta_i, & otherwise \end{cases} \tag{11}$$
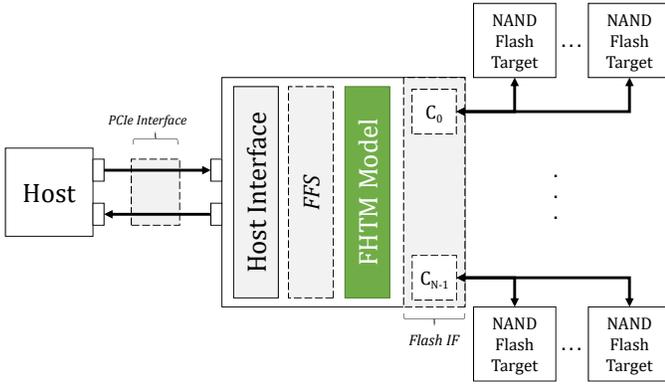
Fig. 2. High-level NVHTM architecture concept. An SSD with a flash file system for address translation. Host interface to manage interoperability with the PCIe interface standard. NVHTM augments the typical SSD structure by incorporating the accelerator into existing data path buses.

where $\beta_i^*$ is the next boost factor value, $\beta_i$, for the $i^{th}$ column; $\widetilde{D}_A$ is the minimum activity duty cycle threshold, and $\beta_{max}$ is the maximum boost factor. The system in (11) is a piece-wise linear function, where $D_A[i]$ is the only quantity that varies between iterations. Within the original CLA algorithmic formulation, $\widetilde{D}_A$ is also specified as a variable quantity. However, as a hardware simplification, this parameter is expected to be extracted from a more complete software model before implementing the model on the hardware platform.

## III. PROPOSED NVHTM ARCHITECTURE

The NVHTM architecture, implemented in VHDL and depicted in Fig. 2, was used to extract area, power, and latency estimates. Operations defined by the mathematical formulation were modeled in hardware as a parallel, pipelined SPU, targeting solid-state drive (SSD) storage devices. Our proposal is that HTM SP logic may be added to the read/write data path, to enable improved parallelism and greater potential for model scalability. Pages of data, containing proximal segment state, are stored within NAND flash target devices through flash channel interfaces. An in-system flash file system (FFS) manages scheduling read and write operations; the host interface exists to translate data from a host CPU to a series of flash commands. Within the hardware model only the NVHTM data path is presented.

Spatial pooling in NVHTM is modeled by augmenting the data path with an overlap pipeline *OVPipe*, inhibition engine *Inheng*, and learning subsystem (*WBPipe* and *WBCntl*). NVHTM auxiliary functions include the channel arbiter (*Charb*), content-addressable memory arbiter (*Camharb*), write-back content addressable memory (*WBCam*), and CAM hit detector (*Camhit*); *WBCntl* is comprised of *Camharb*, *WBCam*, and *Camhit*. NVHTM microarchitecture, depicted in Fig. 3, operates in the following manner:

1) The SSD controller initiates operation by issuing a read command to the flash channel.
2) The overlap pipeline (*OVPipe*) is enabled along with the other functional pipes in the read path. Provides a pointer to the input region vector and configuration data
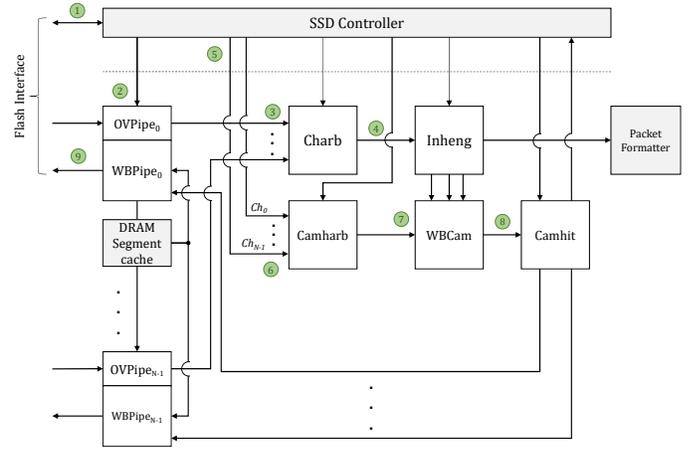


Fig. 3. Microarchitecture for NVHTM SP. Number labels have been applied to indicate the flow of data through SP with learning enabled. Grayed out components indicate that these were not implemented in the NVHTM design but are described to establish context.
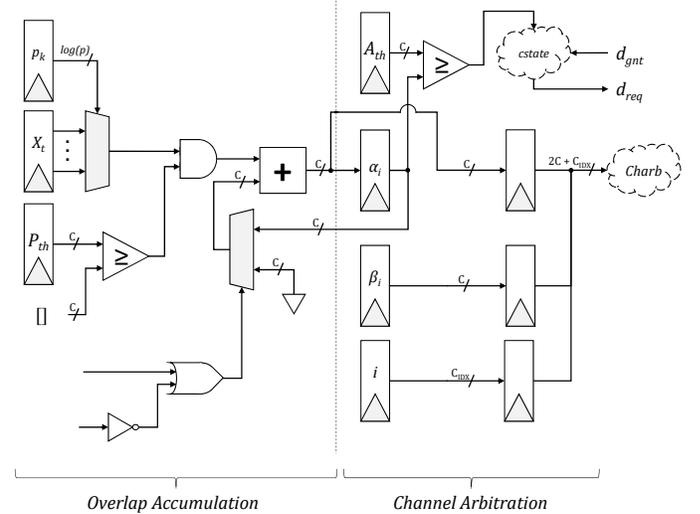


Fig. 4. Overlap pipeline for NVHTM architecture. Configuration data is stored in each of the registers: $P_{th}$ is the permanence threshold; $C_i[k]$ is the permanence for the $k^{th}$ proximal synapse of the $i^{th}$ column; $A_{th}$ is the overlap threshold; $\beta_i$ is the boost constant; $k$ is a pointer used to index into the input vector; $p$ is the length of $X_t$; $p_k$ is the pointer used to index $X_t$; $C$ is the channel width; $C_{IDX}$ is the number of bits used to encode the column index; $psv_0$ is a pipeline stage enable signal; $d_{gnt}$ and $d_{req}$ are the arbitration signals; and $d_{val}$ marks the first valid data on the pipeline. Output of the overlap pipeline is written to a DRAM table and to the inhibition engine. The in-DRAM overlap table is used within the learning phase to update the duty cycle and boost parameters.

broadcast by the host. Data then arrives from the flash interface one beat at a time.
3) Overlap values are selected by the channel arbiter (*Charb*) for boosting and inhibition.
4) After the single columnar overlap and segment state information traverse the pipeline to this point, the inhibition engine (*Inheng*) uses an insertion sort operation to determine the subset of columns that will enter the active state.
5) In the case when learning is enabled, the controller issues column indexes to the content addressable memory arbiter (*Camharb*).
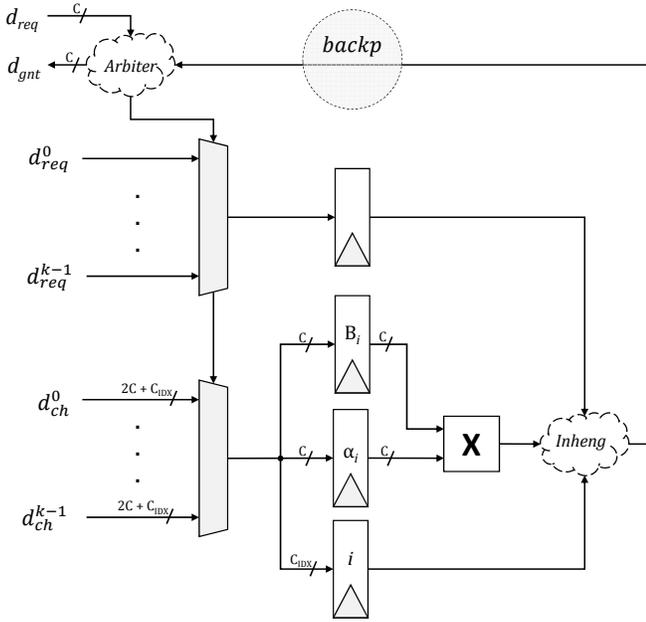
Fig. 5. The channel arbiter (*Charb*) manages transmitting overlap values from multiple overlap engines to the downstream inhibition engine. Furthermore, the overlap values accumulated in *OVPipe* are boosted at this phase prior to being distributed to the inhibition engine. The selection of proximal segment data results in the associated channel receiving a data grant acknowledgment, $d_{gnt}$.
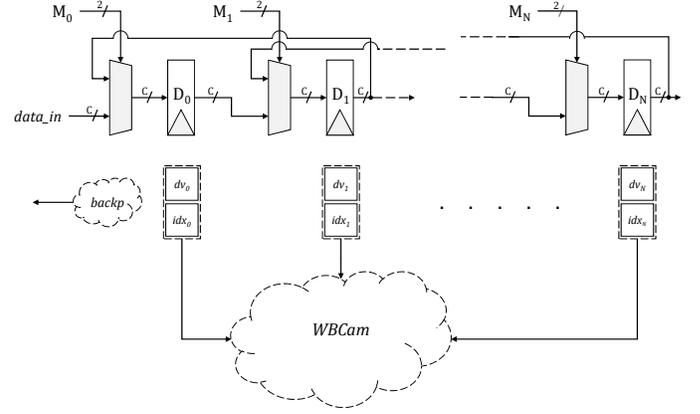


Fig. 6. The inhibition engine for the NVHTM architecture. Control signals are generated by the swap state registers in conjunction with data valid bits and the results of comparison operations. The inhibition engine state is cleared upon the arrival of a new input vector, $X_t$. Note that the column index data path is abbreviated for clarity–it follows a similar path to that of the overlap overlap swap logic, $D_i$.

6) *Camharb* then arbitrates between channels to select a column index to be sent downstream.

7) The write-back content addressable memory (*WBCam*) is used to determine which columns are candidates for learning by comparing the issued column index against valid indexes acquired in the inhibition phase.

8) Hits in the CAM are detected (*Camhit*) and used to control the write-back ALU (*WBPipe*). Timeout signals are forwarded back to the SSD controller to indicate a CAM miss.

9) *WBPipe* uses input from *Camhit* to control how the proximal segment, which is cached in the segment cache, will be written back to memory through the flash interface. State information within the inhibition engine is used to select operands for the ALU within *WBPipe*, which updates each segment one synapse at a time.

After the NVHTM architecture has been configured by the host processor, the SPU waits for an input stimulus, $X_t$, to be broadcast to the SSD. The host processor transmits a packet containing the input vector. Upon reception of the input vector, the FFS generates commands that signal the NAND interface to supply pages of data corresponding to stored proximal segments. Each beat of data on the interface corresponds to synaptic permanences and other SP column states. The codewords of data sourced from the NAND, after low-density parity-check (LDPC decode), are supplied to *OVPipe* one beat of data at a time.

The overlap pipeline, depicted in Fig. 4, is an accumulator that operates synchronous to the flash interface clock. Once the first beat of data arrives at the overlap engine from the LDPC, all subsequent beats of data will become available in the immediately following cycles. Furthermore, because data being read has a predictable structure, a drastically simplified overlap engine was designed. Additional simplification is obtained by exploiting the simplicity of mathematical constructs used by the algorithm.

Latency on the order of tens of cycles are added to the data path due to the vector processor. To put this into perspective, the number of cycles required for the read and write operations are on the order of hundreds to thousands of cycles. Taking into account the interface limitations, additional pipeline complexity is likely to have a negligible impact on performance. As design complexity scales, resulting in larger page sizes, more channels, and higher data rates, this simple model will continue to be feasible and effective. This is because the pipeline has few opportunities for significant critical path limitations: (1) the data interface is only 8-bits to 16-bits wide in most flash implementations, leading to lower computation delays than common 32-bit and 64-bit hardware; (2) there are a small number of inputs to each logical component; and (3) operating at the memory interface frequency provides a large upper bound on allotted cycle time.

As a new $X_t$ input vector arrives from the host, a pointer, $j$, is reset to reference the first bit of the vector. The multiplexer is utilized to select the respective bit within the vector. This first portion of the overlap engine (separated by a dashed line) computes the pre-boosted overlap for the column using an accumulator. Input data arriving from the host is one-to-one matched with the number of potential candidate synapses. Each synaptic value arriving from the flash interface is thresh-olded using a comparator. An AND gate determines that an active-connected synapse was found. In the event that a proximal synapse is both active and connected, the overlap value is incremented. The internal state of the accumulator is reset whenever a new input vector arrives. This complete process takes an amount of time proportional to the page size divided by the channel width–it spans the time required for the read operation plus the overlap pipeline latency.

Following data accumulation in *OVPipe*, the channel arbiter,
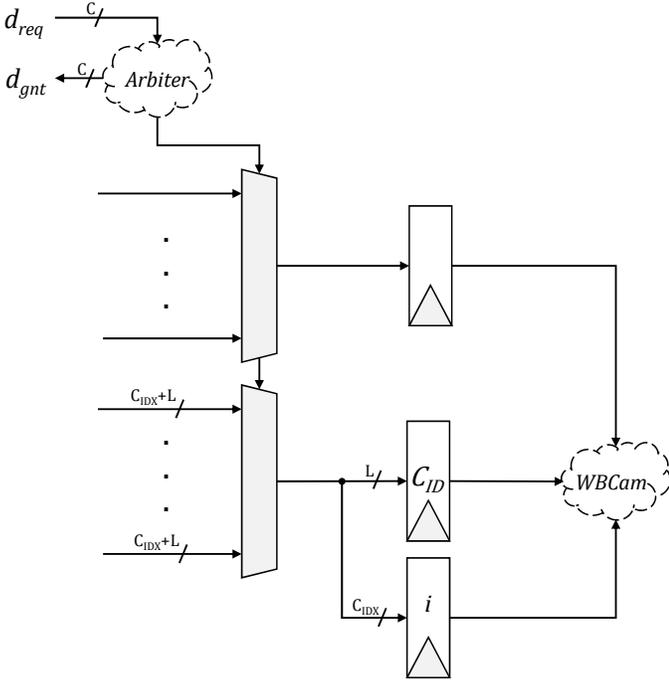
Fig. 7.   The content-addressable memory channel arbiter (*Camharb*) operates similarly to the *Charb*, but does not implement back-pressure because the downstream logic cannot be overrun as with the inhibition engine.



Fig. 8.   The fundamental unit that is chained together to form the *WBCam*. Column index ($cidx_{i-1}$), data valid ($dv_{i-1}$), and channel id ($chid_{i-1}$) from prior stages in the *WBCam* logic are fed into the current stage of the pipeline. Column indexes found in the inhibition queue ($inhidx_i$) are applied to the pipeline and are qualified by valid signals ($inv_i$) to support partially filled inhibition queues.

shown in Fig. 5, is notified that valid data is ready for boosting and inhibition. Only overlap values that have overcome the minimum-overlap constraint issue a valid data request to the downstream logic. A full handshake is computed using a request-grant interface. The purpose of this arbitration is to ensure that only one overlap value is sent downstream and that the other channels are stalled until the inhibition engine is able to service each request in order.

Inhibition for SP is modeled as an insertion sort in which values within the sorted list below a predefined index are discarded. To accomplish this, a pipeline capable of sorting elements in linear time was created. *Inheng*, depicted in Fig. 6, is represented as a multi-directional shifting queue. Overlap values generated from the activation due to the feedforward input are transferred to the inhibition engine after being selected by the channel arbiter. As the overlap pipeline operates on the page data, the inhibition engine attempts to place, and subsequently sort it within the queue structure. Considering the fact that overlap requires hundreds of cycles to compute, the inhibition engine is allotted a large window of time to complete each sort operation. This lead to the design decision to make the sort engine linear, minimizing hardware complexity, whilst still adhering to timing constraints.

This design assumes that the large number of cycles required to compute the overlap eliminates any possibility of overrunning the channel arbiter and in-turn, the inhibition engine. *Inheng* applies back-pressure to *Charb*, stalling it until the prior sort operation has completed. Back pressure acts in a flexible way such that the complete pipe need not be stalled by downstream blockage. Pressure is funneled back upstream stage-by-stage, allowing gaps in the pipe to be filled where
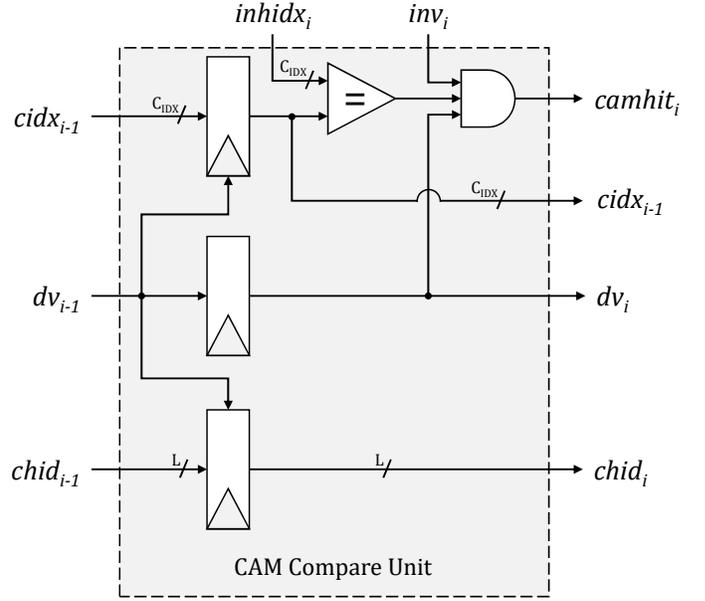
possible. For example, if *Inheng* is applying back-pressure and the output of arbitration is invalid, valid data may still move forward in the pipe.

Data being shifted around in the inhibition queue is $C + C_{IDX}$ bits wide; where $C$ is the flash interface width and $C_{IDX}$ is the number of bits required to encode the column index. The linear sort implemented by the inhibition engine is composed of multiple complementary swap operations governed by sequential control logic. Each swap operation is controlled using a two-bit signal, $M_i$, generated by a swap state register in conjunction with additional combinational control logic. The most significant bit of the swap indicates direction, while the other bit is utilized as a functional enable. Data elements within the queue are swapped using a multiplexer.

Valid bit, column index, and overlap values are moved through the pipeline, controlled by the comparison of overlap data. As new data is loaded, swap requests propagate through the queue until a valid swap operation is no longer possible. The swap logic checks for validity between the two elements and compares their magnitude, taking into account queue fullness. *Inheng* is deactivated after completing the swap operation as it waits for new overlap data to become available.

The control signals are specified to allow all registers to shift in the same direction when new data is loaded, and to subsequently sort that data. After all HTM columns have been processed by *Inheng*, the column indexes corresponding to the overlap values stored within the queue are used to generate an SP output packet due for transmission to the host processor. Column indexes generated by the inhibition phase are also used in the learning phase to influence proximal segment update.

Read data for SP flows from *OVPipe*, to the *Charb*, finally

ending in *Inheng*, which supplies output to a packet formatter that emits through the host interface. Learning follows a slightly different data flow and must immediately follow the completion of inhibition for all columns. Indexes for all columns in *Inheng* are fed into the *WBCam* for use in the learning phase. Proximal segment update is initiated by the SSD controller, which must schedule reads for all columns in the network and also issue corresponding notifications to the content-addressable memory arbiter (*Camharb*), shown in Fig. 7. The SSD controller is notified of hits and misses in the CAM by *Camhit*.

*WBCam* is comprised of chained compare units, depicted in Fig. 8. For each stage other than the first, if a hit is detected, the data is invalidated in the next stage in the pipeline to avoid another subsequent comparison. This feature is requisite to avoid invalid detection of timeouts in the pipeline. For all compare components other than the first ($i > 0$), the data valid output is defined by (12).

$$dv_i = dv_{i-1} \wedge \overline{inhv_{i-1} \wedge (cidx_{i-1} = inhidx_{i-1})} \quad (12)$$

Proximal segment data is read from the memory, and subsequently cached in the lower latency DRAM. Upon completion of the read operation, the data is written from the segment cache in DRAM back to the SSD NAND array through *WBPipe*. Operands for the write-back operation are selected based upon whether a hit was detected in *WBCam*, and upon the overlap duty cycle. Proximal segment state information is updated in the *WBPipe* as well. In this scheme, the memory can queue up several reads to occur in order, followed by segment updates.

Hits in the CAM are associated with their originating channel via the *Camhit* component, which utilizes a series of equality comparators to aggregate hits corresponding to a channel of interest. This scheme also follows the assumption that each channel will only issue a single column index into *WBCam* at a time to remove the possibility for ambiguity with respect to detected hits. An additional comparison is made with the final phase of the *WBCam* and this index is assigned to the timeout operation logic. Column indexes are indicated as having timed out if they have been compared against all data in the inhibition engine without having matched any other data. A hit detector unit for a single channel is depicted in Fig. 9. In an SSD with $N$ channels, this logic is repeated $N$ times.

Using the results of CAM hits and misses, the SSD controller modifies the input operands to the *WBPipe* ALU. A hit in the CAM indicates that the column was active and therefore a valid candidate for learning. All segments are still re-read and written back to memory in the event that an update is pending for the column. Segments may be updated due to being in the active state following feedforward input or as a result duty cycle updates. In the case of an active duty cycle, the whole segment need not be updated, only the active duty cycle and boost factor. Columnar activity and overlap duty cycle result in an update being made to the entire set of synapses on the segment. Duty cycle parameters cached in the DRAM as proximal segments are first read from memory by the *OVPipe* to control the scheduling of segment updates.
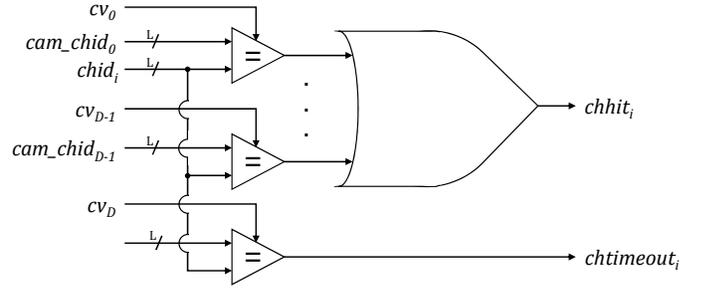


Fig. 9. The *Camhit* component is comprised of comparators used to determine whether the ID for the $i^{th}$ channel, $chid_i$, has been matched in the CAM. It returns a notification to the SSD controller indicating the event. Valid signals, $cv_{0..D}$, are used to qualify each comparison. A timeout signal, $chtimeout_i$, may also be asserted indicating that a compare operation timed out, resulting in a CAM miss.

State information regarding the proximal segment is updated by the write-back data pipeline shown in Fig. 10. This pipeline is capable of updating segment duty cycle and boost parameters in addition to the synaptic permanences on the segment itself. Flash interface data is redirected to the pipeline associated with control information provided by the SSD controller. With this topology, data can be streamed in and out of the update pipe at the flash interface speed to reconstruct the original proximal segment page, given the correct control sequence. Output from *WBPipe* may be sourced from the duty cycle pipeline, the old boost factor value register, or the proximal segment update data path.

Configuration inputs, labeled $y_{1..9}$, each correspond to various pre-calculated constants necessary to compute updates in the learning phase. Each constant is defined in (13). The segment update formulas were modified to be more amenable to the NVHTM design concept, meaning that variable parameters were separated out and divisions were represented as pre-calculated constants instead of in the hardware. This modification made it possible to eliminate the need for dividers and also reduced the number of multipliers to one per write-back pipeline. Although, this could be improved further by using an out-of-order execution core, capable of issuing commands from a fixed number of agents to a fixed set of resources. Wavefront allocators have been employed in the literature to provide effective resource management in hardware architectures.

$$y_1 = (\tau_D - 1)/\tau_D \quad (13a)$$

$$y_2 = 1/\tau_D \quad (13b)$$

$$y_3 = (1/\widetilde{D}_A)(1 - \beta_{max}) \quad (13c)$$

$$y_4 = \beta_{max} \quad (13d)$$

$$y_5 = P_{th}/10 \quad (13e)$$

$$y_6 = P_{inc} \quad (13f)$$

$$y_7 = P_{dec} \quad (13g)$$

$$y_8 = P_{inc} + P_{th}/10 \quad (13h)$$
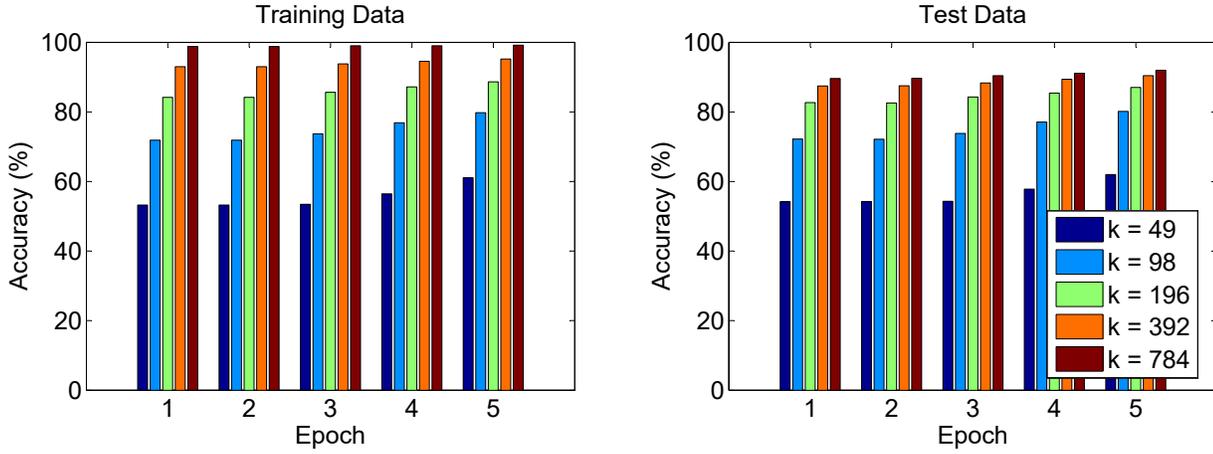
$$y_9 = P_{dec} + P_{th}/10 \quad (13i)$$

Fig. 11. Classification results for the SP math model, with varying network sizes, $K$. Network performance attained near optimality after the first epoch, with nominal improvement through subsequent iterations.
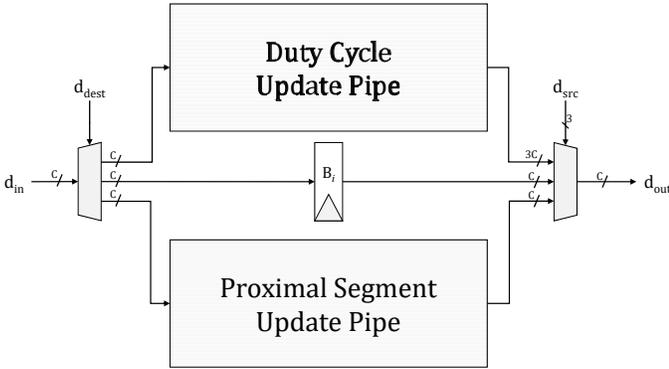


Fig. 10. High-level view for the write-back pipeline, *WBPipe*. This structure enables data from the flash interface to be steered toward either the duty cycle or proximal segment update pipe. Furthermore, an additional staging register is added to forward the original boost factor during its update phase to allow for the new boost value to be rejected in the event that an active duty cycle is above threshold. Sequencing is controlled by the SSD controller, allowing data to be sourced from either of the aforementioned pipes.
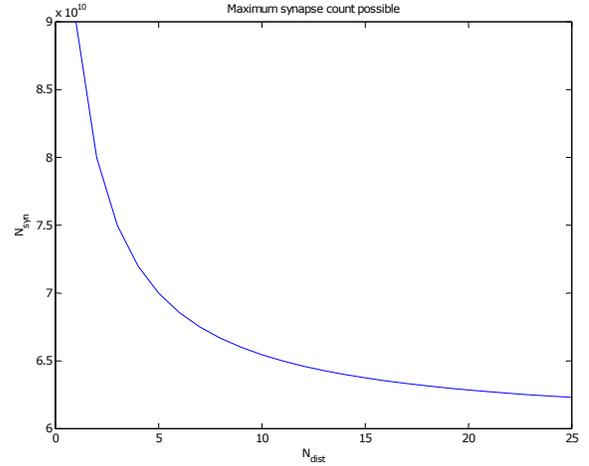


Fig. 12. Maximum number of synapses possible for HTM network, as a function of segment size, assuming a fixed memory size ($Mem_{sz} = 240GB$), and fixed page size ($P_{sz} = 32kB$); channel width assumed to be 16-bits; column index width, $C_{IDX}$ is assumed to be 24 bits–the number of cells that may be accessed is on the order of millions.

## IV. RESULTS

### A. Classification

Results for SP were obtained by simulating the mathematical model using a single thread process on a fixed platform. Execution time required to simulate five epochs across five unique, randomly initialized networks grew from approximately two hours ($K = 49$) to about a week ($K = 784$); this time included training, test, and data collection times for the spatial pooling algorithm. The Matlab model was not optimized for speed, instead, is a demonstration of the functional fidelity of the model.

Classification performance, shown in Fig. 11, was found to reach near optimality after the first epoch, following the law of diminishing returns for subsequent epochs. Accuracies of 99.17% for training and 91.98% testing were observed for $K = 784$. The impact of training for additional epochs appeared to play a more significant role in networks with fewer columns. This suggests that an SP-SVM (support vector machine) configuration may be trained with a single SP training epoch for larger network sizes, while still achieving comparable performance. There also seems to be an upper limit on the benefits of network scaling; the improvement from 392 to 784 columns was significantly smaller than that observed between 49 and 98 columns.

### B. Network Model Scalability

NVHTM has the potential to represent large scale networks constrained by the memory size and network configuration. Taking model configurations into consideration, a function was derived to ascertain the potential network size. Total memory and page sizes were used as key parameters. The number of distal segments was considered to demonstrate how this model could be expanded for TM computation support. Maximum network size is limited to the number of pages, because each page represents a single proximal or distal segment. The number of pages within the memory can be simply defined as

$$N_{pages} = \frac{Mem_{sz}}{Page_{sz}} \qquad (14)$$

where $Mem_{sz}$ is the size of user memory space in bytes and $Page_{sz}$ is the size of pages in bytes. Each column has a set of distal segments associated with it, all of which are connected to cells within the respective column. Furthermore, columns are assumed to only have one proximal segment assigned to them. This is used to determine the number of proximal segments as a function of memory size using the relationship

$$N_{prox} = \frac{N_{pages}}{1 + N_{dist}} \qquad (15)$$

where $N_{prox}$ is the number of proximal segments; and $N_{dist}$ are the number of distal segments per column: equivalent to the cells per column multiplied by the number of distal segments per cell. Each column is assumed to have a single proximal segment.

Using these equations, the quantity of synapses may also be defined by noting that each proximal segment page has a fixed number of synapses three words less than the page size: the first three words of data contain the overlap and active duty cycles, in addition to the boost factor. Distal segments, assumed to be represented in a sparse format, may only represent a number of synapses less than half the page size (in words), because synapse values are stored in index-value pairs [18]. This is defined using the following system

$$N_{syn} = N_{prox} \frac{P_{sz} - 3C}{C} + N_{prox} N_{dist} \frac{P_{sz} - C}{C + C_{idx}} \qquad (16)$$

where $N_{syn}$ is the number of synapses within the model; $P_{sz}$ is page size in bytes (flash interface width); $C$ is the number of bytes required to represent a synaptic permanence, assumed to be the flash data interface width (8bits/16bits); $C_{idx}$ is the number of bytes required to represent the cell index for the sparse distal segments. A sample plot of (16) is shown in Fig. 12.

### C. Power, Area & Latency

Hierarchical synthesis was explored for *WBPipe* and *WBCntl* to make synthesizing these larger designs more feasible. Furthermore, instantiated components were compacted to reduce area and an iterative workflow was used to ensure that DRC and LVS requirements were met. More aggressive routing strategies such as wrong-direction routing, ripping and congestion analysis, were required to ensure that all of the overflows between nets would be properly routed. Furthermore, some primary I/O pins had to be manually moved to ensure routability. *WBCntl*: composed of *Camharb*, *WBCam*, and *Camhit*: was created to explore block-based hierarchical synthesis. The power results for each of these components were obtained by extracting the parasitic netlist of *WBCntl*.

Area footprints were generated by Mentor Graphics Pyxis IC design suite as summarized in Fig. 13. The dimensions of each component are provided, with a total footprint of $30.538mm^2$. However, the estimated area required for an actual implementation must take into account the number of
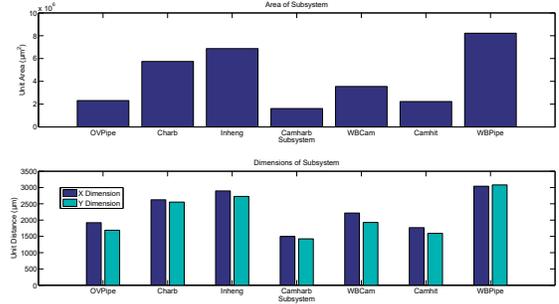


Fig. 13. Area (top) and dimensions(bottom) for each of the synthesized RTL components. The read and write pipe logic are expected to consume the greatest amount of area, because they are duplicated per channel, resulting in an 8x-20x increase in area for these components.
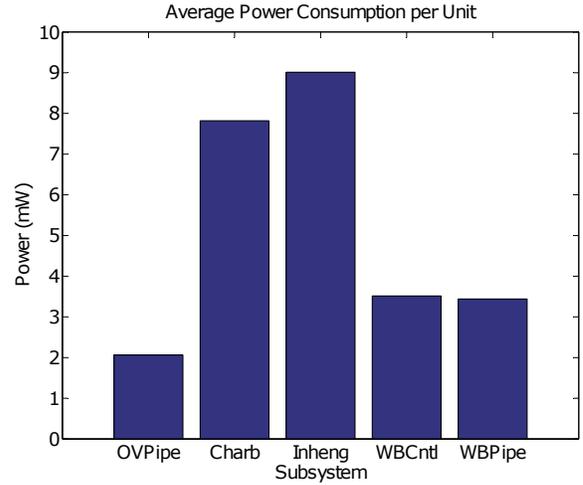


Fig. 14. Average power consumed for each of the synthesized RTL components.

times that each component is expected to be instanced in an actual SSD NVHTM. Estimated area is defined as

$$A_{NVHTM} = N_{ch}(x_0 + x_1) + x_2 + x_3 + x_4 + x_5 + x_6 \qquad (17)$$

where $N_{ch}$ is the number of channels, $x_0$ and $x_1$ are the areas for *OVPipe* and *WBPipe*, respectively; the other $x_i$ terms correspond to each of the other components. Under this assumption, the area becomes $104.26mm^2$.

A parasitic extraction netlist was derived from the layout and evaluated using the Eldo analog simulator. Hardware created at this phase was found to consume power on the order of a few milliwatts, as depicted in Fig. 14. Components consuming the most power did so due to having a larger number of memory elements switching state at each clock cycle. Average power consumed per component was $5.171mW$–the estimate for an 8-channel SP model is $64.394mW$.

## V. CONCLUSION & FUTURE WORK

As our understanding of the underlying mechanisms that drive intelligence continues to expand, applications that leverage this body of knowledge will increase. It is likely that

machine intelligence, applied to enterprise-scale big data analysis of unstructured data, will serve as catalyst for future innovations. Advances in algorithms will, in turn, lead to the creation of new architectures. Storage processors present a promising opportunity for product differentiation in this space. Companies may provide robust virtual platforms capable of handling pattern recognition workloads at scale with low cost.

Through this exploration, a conceptual design and analysis have been presented for an NVHTM spatial pooler. The impact of augmenting a storage unit with processing capabilities degrades the upfront latency for reading data, consumes additional area resources, and may potentially limit the usefulness of this storage unit for other applications. Despite this, there are several benefits to taking the SPU approach:

1) Parallelism can be more thoroughly exploited than if the algorithm were implemented by a multi-core processor
2) Vital resources within the host are conserved, allowing the storage unit to act as an accelerator at scale
3) A key challenge for storage elements is I/O availability: an SPU does not require any additional I/O
4) Memory-mapped configuration facilitates the design of a microarchitecture that is simple to implement, configure, and extend
5) An in-path SPU may be sold as locked IP that may be dropped into other SoC environments
6) Significant power savings are obtained over in-host processing, which operates at an order of magnitude higher frequency
7) Scalability of the design is significantly improved over external accelerator designs, which are bottlenecked by memory bandwidth limitations
8) This design can be scaled to other emerging memory devices that offer similar competitiveness as flash, such as memristors or PCM.

A clear explanation of the model utilized to implement the HTM spatial pooler has been presented, providing insight into NVHTM design limitations and benefits. SP, originally developed with high-level programming in mind features significant branching behavior and unclear parameter limits. Despite these challenges, a simplified hardware model has been presented along with an accompanying system model. Power, area, and latency estimates were extracted from each phase of the design process to acquire a baseline for feasibility analysis. NVHTM, paired with SVM for classification, present results comparable to those found within the literature for HTM MNIST (91.98%).

In conclusion, the large delays characteristic of SSDs (when compared to main memory or cache) mean that the latency added by the NVHTM pipeline is orders of magnitude less than the baseline latency for a standard SSD. This supports the case for employing storage processor units in hardware, as has been discussed in the literature. However, the specialization required by the hardware limits the potential for deploying this to large scale markets. Implementing the SPU functionality on a reconfigurable platform integrated into the SSD would be an attractive option, because it has the potential to improve design scalability. This framework can be extended to the temporal memory of the HTM, which is designed to model the inter-

pattern association across various time steps. However the challenge to modeling it is the unbounded segment growth. To manage this specialized techniques such as dynamic pruning and time multiplexing have to be employed.

## REFERENCES

[1] J. Hawkins and S. Blakeslee, *On Intelligence*. Times Books, 2004.

[2] D. George and J. Hawkins, "A hierarchical bayesian model of invariant pattern recognition in the visual cortex," in *In Proceedings of the International Joint Conference on Neural Networks. IEEE*, 2005, pp. 1812–1817.

[3] D. 'George and J. Hawkins, "Towards a mathematical theory of cortical micro-circuits," *PLoS Comput Biology*, vol. 5, no. 10, 10 2009.

[4] N. Inc., "Hierarchical temporal memory: Concepts, theory, and terminology," Numenta Inc., 791 Middlefield Road Redwood City, CA 94063, Tech. Rep., 3 2006.

[5] D. George and B. Jaros, "The htm learning algorithms," Numenta Inc., 791 Middlefield Road Redwood City, CA 94063, Tech. Rep., 3 2007.

[6] N. Inc., "Hierarchical temporal memory including htm coritcal learning algorithms," Numenta Inc., 791 Middlefield Road Redwood City, CA 94063, Tech. Rep. 0.2.1, 9 2011.

[7] J. Mnatzaganian, E. Fokoue, and D. Kudithipudi, "A mathematical formalization of hierarchical temporal memory's spatial pooler," in *IEEE Transactions on Neural Networks and Learning Systems*, 1 2016, pp. 1–11.

[8] Y. Cui, C. Surpur, S. Ahmad, and J. Hawkins, "Continuous online sequence learning with an unsupervised neural network model," *CoRR*, vol. abs/1512.05463, 2015.

[9] Numenta, "The science of anomaly detection," 3 2015.

[10] S. Ahmad and J. Hawkins, "Properties of sparse distributed representations and their application to hierarchical temporal memory," 3 2015.

[11] X. Zhou and Y. Luo, "Implementation of hierarchical temporal memory on a many-core architecture," Master's thesis, Halmstad University, PO Box 823 SE-301 18 Halmstad, Sweden, 12 2012.

[12] R. W. Price, "Hierarchical temporal memory cortical learning algorithm for pattern recognition on multi-core architectures," Master's thesis, Portland State University, 1825 SW Broadway, Portland, OR 97201, 1 2011.

[13] V. Chandola, A. Banerjee, and V. Kumar, "Anomaly detection: A survey," *ACM Comput. Surv.*, vol. 41, no. 3, pp. 15:1–15:58, Jul. 2009. [Online]. Available: http://doi.acm.org/10.1145/1541880.1541882

[14] D. Fan, M. Sharad, A. Sengupta, and K. Roy, "Hierarchical temporal memory based on spin-neurons and resistive memory for energy-efficient brain-inspired computing," *Computing Research Repository*, 1 2014.

[15] X. Chen, W. Wang, and W. Li, "An overview of hierarchical temporal memory: A new neocortex algorithm," in *Modelling, Identification Control (ICMIC), 2012 Proceedings of International Conference on*, 6 2012, pp. 1004–1010.

[16] A. Liddiard, J. Tapson, and R. Verrinder, "A robust implementation of the spatial pooler within the theory of hierarchical temporal memory (htm)," in *Robotics and Mechatronics Conference (RobMech), 2013 6th*, 10 2013, pp. 70–73.

[17] P. Vyas and M. Zaveri, "Verilog implementation of a node of hierarchical temporal memory," *Asian Journal of Computer Science and Information Technology*, no. 3, pp. 103–108, 7 2013.

[18] A. M. Zyarah, "Design and analysis of a reconfigurable hierarchical temporal memory architecture," Master's thesis, Rochester Institute of Technology, 1 Lomb Memorial Dr, Rochester, NY 14623, 6 2015.

[19] V. Puente and J. A. Gregorio, "Claasic: a cortex-inspired hardware accelerator," 4 2016, pp. 1–13.

[20] M. Prezioso, F. Merrikh-Bayat, B. Hoskins, G. Adam, K. K. Likharev, and D. B. Strukov, "Training and operation of an integrated neuromorphic network based on metal-oxide memristors," *Nature*, vol. 521, no. 7550, pp. 61–64, 2015.

[21] S. Hu, Y. Liu, Z. Liu, T. Chen, J. Wang, Q. Yu, L. Deng, Y. Yin, and S. Hosaka, "Associative memory realized by a reconfigurable memristive hopfield neural network," *Nature communications*, vol. 6, 2015.

[22] G. W. Burr, R. M. Shelby, S. Sidler, C. Di Nolfo, J. Jang, I. Boybat, R. S. Shenoy, P. Narayanan, K. Virwani, E. U. Giacometti *et al.*, "Experimental demonstration and tolerancing of a large-scale neural network (165 000 synapses) using phase-change memory as the synaptic weight element," *IEEE Transactions on Electron Devices*, vol. 62, no. 11, pp. 3498–3507, 2015.

[23] J. Hsu, "Ibm's new brain [news]," *IEEE Spectrum*, vol. 51, no. 10, pp.

[24] B. V. Benjamin, P. Gao, E. McQuinn, S. Choudhary, A. R. Chandrasekaran, J.-M. Bussat, R. Alvarez-Icaza, J. V. Arthur, P. A. Merolla, and K. Boahen, "Neurogrid: A mixed-analog-digital multichip system for large-scale neural simulations," *Proceedings of the IEEE*, vol. 102, no. 5, pp. 699–716, 2014.

[25] H. B. Project. (2014) Brainscales project. [Online]. Available: https://www.technologyreview.com/s/536326/ibm-tests-mobile-computing-pioneers-controversial-brain-algorithms/

17–19, 2014.