

# A Nonvolatile Approximate Lookup Table Optimized for GPU Acceleration

Snehasis Dey,

Assistant Professor,

Dept.of Electronics & Telecommunication Engineering  
College of Engineering Bhubaneswar

**Abstract**—In this letter, we design a nonvolatile approximate lookup table, named NVALT, to greatly accelerate general public utilities (GPUs) computation. Our architecture maintains high frequency input patterns within an approximation NVALT to simulate each application's functionality. To provide an approximative output, NVALT looks for and returns the stored data that most closely matches the input data. By taking use of the analog characteristics of the nonvolatile content addressable memory, we define a similarity measure that is suitable for binary representation. To regulate the precision required for user requirements, our approach manages the ratio of the application executing between the approximate NVALT and correct GPU cores. Our tests on seven common GPU applications demonstrate that, on average, NVALT may increase energy computation by 4.5 and performance by 5.7 while offering an average relative inaccuracy of less than 10% when compared to the baseline GPU.

## I. INTRODUCTION

INTERNET of Things (IoT) increases the number of smart devices and the rate of data generation around the world [1], creating a significant demand for high speed, and efficient parallel processors, such as general public utilities (GPU). However, GPU cannot learn existing patterns in workload and adaptively process the data [2]. When repeatedly running a single application, e.g., fast Fourier transform (FFT), on a GPU for different input data, the GPU performs a set of computations repeatedly without learning the functionality. By performing nonconventional brain-like computation, e.g., neural network and neuromorphic computing solutions perform, the cores avoid costly repeated computations by learning sets of functionality and approximately model them [3]. This technique is effective for applications which can accept a level of approximation. Many real world algorithms, such as machine learning, are statistical in nature and will accept some inaccuracy in their computation. Several image and video processing applications accept error in computation without losing the precision required by user [2]. We add an efficient brain-like processing unit beside the GPU cores to accelerate the computation with approximation techniques.

Prior work tried to improve efficiency and accelerate the GPU/CPU computation by enabling approximation, such as voltage over-scaling [4], precision scaling [5], designing approximate circuits [6], [7], and memoization using resistive acceleration [8]–[12]. Work in [3] accelerates programs by utilizing a neural network placed beside the GPU cores to produce

approximate results. Although in CPU this technique provides a large energy/performance improvement [14], the advantage of neurally based approximation in GPU is minor. Several other designs use the high density and zero leakage power of nonvolatile memories for memoization or enabling approximation in GPU and CPU cores [8], [10], [13]. Work in [8] enabled GPU approximation by using content addressable memories (CAMs) beside floating point units (FPUs) and approximately retrieving data at run-time. This is a promising technique to save the GPU energy, but it cannot improve the performance, and also does not have a large impact on the overall GPU energy consumption, considering data movement and other non-FPU units.

In contrast, we propose a hardware approximation technique that uses resistive content addressable memories to accelerate the GPU computation. We exploit the analog characteristics of the nonvolatile approximate lookup table, called NVALT, with nearest distance search capability. Our design learns and models the functionality of different applications by storing high frequency patterns for each application on an NVALT. At run-time, instead of processing data on the inefficient GPU core, our design searches NVALT to find the data most similar to the input operands. Our similarity metric considers the impact of each bit index to provide the best match. Our experimental evaluation on advanced micro device (AMD) Southern Island GPU architecture, running seven different applications, shows the enhanced GPU, as compared to the baseline GPU, can achieve 5.7 energy and 4.5 performance improvement, while providing less than 10% quality loss.

## II. PROPOSED NVALT

### A. NVALT Program Acceleration

GPU workload exhibits significant data similarity and locality. This locality increases as we go through the IoT domain with larger datasets. Several basic programs, such as FFT or image processing applications, consist of repeated building blocks, where each has many addition and multiplication operations. Running these algorithms on the conventional processors (e.g., CPU or GPU) is slow and requires large energy consumption. These applications can be approximately processed on the NVALT block. We exploit data locality to model

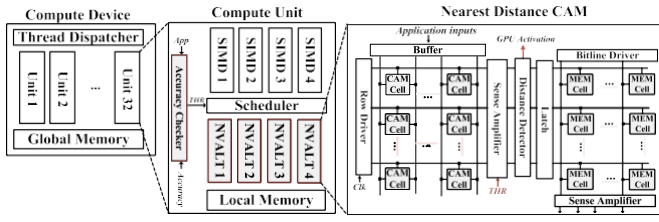


Fig. 1. Overview of the enhanced GPU with approximate NVALT block.

several basic functions on the GPU with a fast and efficient approximate computing unit.

Fig. 1 shows the overview of the structure of AMD's Southern Islands GPU, enhanced with NVALT accelerators. NVALT blocks are placed along side of each single instruction multiple data (SIMD) lane. When the approximate application is launched on GPU, the scheduler block runs the application on an NVALT block instead of sending them to GPU core to process. An NVALT block approximately processes these basic functions inside a lookup table instead of precisely processing them. NVALT uses offline profiling for each program to find and store common input patterns and the corresponding outputs. At runtime, our design searches the input data stored in the NVALT and returns the most similar pattern as an approximate output. The accuracy of depends two main metrics.

- 1) *Similarity Metric*: In traditional cores the numbers represent by the fixed-point or floating point binary values, where the definition of similarity can be different based on representation. For instance, our prior work [10] considers Hamming distance as a metric to find the most similar row in a lookup table. This metric cannot provide high computation accuracy because it does not consider the impact each bit index has on the computation. Designing a lookup table/CAM which can find the exact nearest value requires large and inefficient peripheral circuitry. In this letter we propose a simple technique which changes the weight of each bit index on the CAM structure and considers the binary weights on the search operation. We detail the implementation in Section II-B.
- 2) *NVALT Size*: The lookup table size directly impacts the computation accuracy. A large size NVALT can store a large number of patterns and increase the chance of a close match. Larger lookup tables result in the increase energy and degraded performance.

NVALT works for the applications which accept approximation in their computation. Applications must have limited number of input/output signals to efficiently fit within the NVALT table. We model the computation of these applications by storing their highly frequent input/output patterns inside a table. Profiling mode examines of the input data to find the most common occurrences.

Our evaluation shows that running all applications on the approximate NVALT may not result in high energy saving or speedup. We need to increase the size of the NVALT to store a sufficient number of patterns to ensure the acceptable computation accuracy. The increased size reduces the energy and performance efficiency that NVALT achieves during the search operation. The low energy and performance advantage of neural network-based GPU acceleration in [3] comes from using large and inefficient neural network blocks to provide

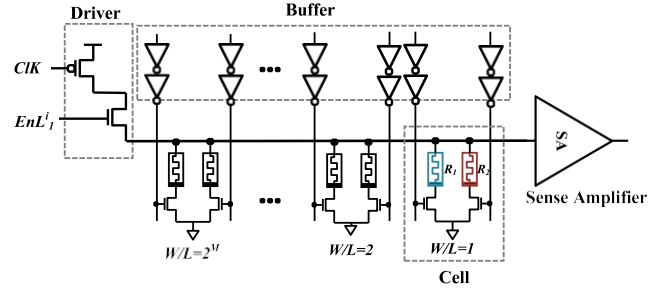


Fig. 2. Structure of CAM in NVALT block capable of searching for nearest distance row.

Data	Conv. Cell		NVALT Cell		Search Operation	
	R1	R2	R1	R2	SL/SL'	ML
0	L	H	H	L	Vdd/0	Precharge
1	H	L	L	H	0/Vdd	

Fig. 2. Structure of CAM in NVALT block capable of searching for nearest distance row.

data runs on the accurate GPU core, while the rest is run on the approximate NVALT. Our design calculates the distance of the input data with the stored value on the NVALT, and if the distance is larger than a threshold ( $>THR$ ), it dynamically assigns this input data to an accurate core. The value of  $THR$  is determined based on the application type and user accuracy requirement.

## B. NVALT Hardware Support

A conventional CAM consists of a CAM cell array, row driver, input buffer, and sense amplifier. Before each search operation, the match lines (ML) of all CAM rows are precharged to  $V_{dd}$ . During the search operation, input data is distributed to all CAM rows using an input buffer. This buffer strengthens the input signals to ensure all rows receive the input signals at a similar time. During the search operation, the MLs in all CAM rows are discharged as long as at least one bit difference exists. We use the timing characteristic of ML discharging current to differentiate rows with different number of mismatches. In conventional CAM, any cell with mismatch discharges the ML. A larger number of mismatches results in a higher ML voltage drop, which can be detected by the sense amplifier. To detect a row with minimum mismatch, i.e., closest Hamming distance row, we need to find the row which discharges last. The sense amplifier tracks the ML voltages in all rows, until it finds the slowest discharging one. The design is complicated and requires additional circuitry, such as counters, while also taking a long time to determine the best matched row.

To address this issue, we use inverse CAM, proposed in [10], to design a CAM with nearest distance search capability. The CAM cells in our proposed NVALT work inversely, compare to the conventional CAM. Fig. 2 shows the functionality of NVALT cells storing inverse resistance values in match and mismatch cases. NVALT cell discharges the ML in case of matching data to the stored values, while a mismatch ML stays charged. A row which has more matched bits, i.e., a closest Hamming distance row, creates a faster ML discharging current than other rows. We detect the nearest distance row by finding the first row to discharge the ML.

Fig. 2 shows a CAM in single NVALT stage which is capable of searching for the row with the closest value to the provided inputs. The proposed NVALT consists of an input buffer, a



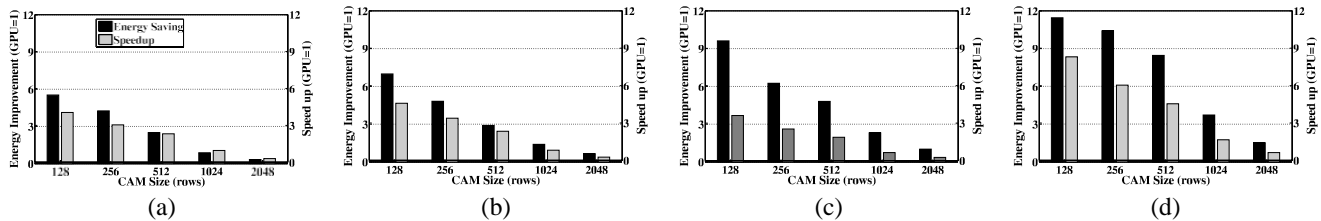


Fig.3. Energy improvement and speedup of the NVALT in different sizes. (a) Black Scholes. (b) FFT. (c) Sobel. (d) inversek2j.

row driver to selectively precharge ML of each row, and a sense amplifier to detect the first discharged row. Unlike other designs which use the hamming distance criteria, our design considers the impact of each bit index on the search operation of the NVALT block. We exploit different access transistor sizes for different bit indices. Based on the binary weight of an unsigned integer value, each cell in position  $i$ th has access transistors which are 2 larger than the cell in the  $i$ th adjacent bit. This results in 2 higher ML discharging current in each match cell compared to its adjacent least significant bits. The asymmetric access transistors weight each bit index on the search operation to find the closest row.

### III. EXPERIMENTAL RESULTS

#### A. Experimental Setup

We evaluate the efficiency of the proposed NVALT on the AMD Southern Island GPU, Radeon HD 7970 device, which is a recent GP-GPU architectures. We modified the source code of Multi2sim, a cycle-accurate simulator [15], to integrate the NVALT within GPU. We used McPAT tool [16], to measure the energy consumption of GP GPU and other change on the GPU architecture, including the registers and first-input, first-output (4KB/SIMD). We also used CACTI to measure the energy consumption of the memory keeping the highly frequency patterns corresponding to each application. To measure the power consumption of the NVALT, we use HSPICE circuit level simulation in 45nm TSMC technology. We use VT

EAM memristor model [17] for our memory design simulation with  $R_{ON}$  and  $R_{OFF}$  of 10k  $\Delta$  and 10M  $\Delta$ , respectively. In terms of reliability, the endurance of resistive memory limits to  $10^7$  write operations, our design addresses this issue by limiting the number of writes to once for each application. We compare the efficiency of proposed design by running seven general OpenCL applications: *Black Scholes*, *FFT*, *Sobel*, *inversek2j*, *Laplacian*, *Convolution*, and *Binarization*.

#### B. NVALT Accuracy Tuning

NVALT can be used next to GPU SIMD lanes to approximately process different applications. We use the approximate building block as a stand-alone computing unit to approximately process the desired applications. In this mode, the NVALT accuracy tunes using the size of the table that we are using. As we explained in Section II, there is an efficiency and accuracy tradeoff in choosing the best NVALT size to process the data. Small size table can provide significantly higher energy savings and performance improvement at the cost of decreased accuracy. Increasing the table size increases the computation accuracy by storing more common patterns, thus increasing the possibility of close hit. However, NVALT with many rows requires larger interconnects and column driver to distribute the data simultaneously among all CAM rows. TABLE I

TUNING ACCURACY IN ROW-TUNABLE NVALT BLOCK

# of rows	128	256	512	1024	2048
<i>Black Scholes</i>	21%	13.1%	7.2%	4.0%	1.2%
<i>FFT</i>	8.3%	5.9%	3.1%	1.9%	0.7%
<i>Sobel</i>	13.4%	6.4%	4.2%	2.1%	1.4%
<i>Inversek2j</i>	17.3%	12.0%	8.3%	5.1%	2.9%

This energy overhead takes the search energy (the energy consumed by the array), for CAM larger than 1024 rows.

Our design is able to dynamically change the number of active rows on the NVALT (128-row granularity) in order to get the maximum accuracy for each application. We enable this functionality by chaining the structure of the row driver. Fig. 3 shows the energy and performance improvement of four general applications on the proposed enhanced hardware. For each CAM size, the energy and performance have been normalized to the conventional GPU not using NVALT. Table I shows the computation accuracy of each application for different CAM sizes. Our results show the applications have different accuracy sensitivity to the NVALT size, as each application

requires a different number of rows to provide acceptable quality of service. To achieve less than 10% quality loss, our design needs to use an NVALT up to 512 rows. The result shows that our design achieves 5.3 energy improvement and 3.2 speedup, when selecting the optimal NVALT size for each application. The main issue of the row tuning technique is its weakness in providing energy/performance advantage for high quality computation. To provide less than 4% quality loss, our design does not save energy or performance.

#### C. GPU-NVALT Computing

To address the quality issue in row tuning technique, we need to understand why very large NVALT blocks still result in sizable error. NVALT shows high error running inputs

that have low similarity to the stored values. Although the percentage of data with low similarity makes up a small fraction of each workload, the impact of them on accuracy is high. In most cases, these infrequent patterns of low similarity data are the main information of the input. For example, in image processing, the image edges have the low frequency compared to the image background, but the edges contain the main information of image. To improve the NVALT accuracy, our

design dynamically finds the far input values and assigning them to accurate SIMD cores to process. The ratio of running data on the approximate and precise core determines the level of accuracy that our design can achieve. Fig. 4 shows the maximum energy improvement and speedup that our design can achieve in each NVALT size when the threshold value changes. The THR has been set in order to ensure that GPU quality loss is

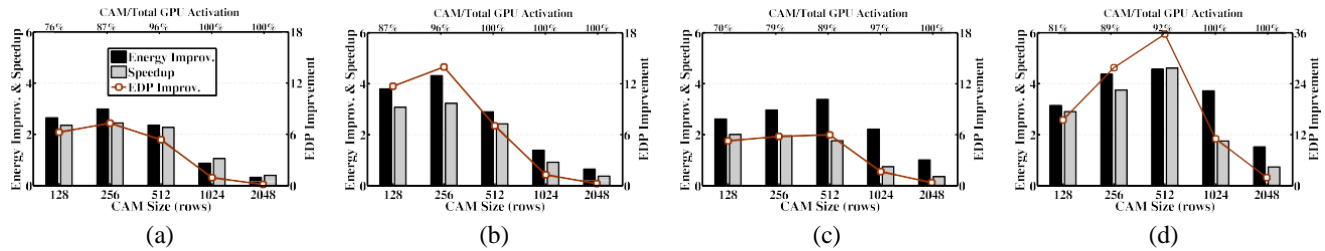


Fig.4. Energy improvement, speedup, and EDP of enhanced-GPU with NVALT in different CAM sizes while ensuring less than 4% quality loss. (a) BlackScholes. (b) FFT. (c) Sobel. (d) inversed2j.

TABLE II  
 EDP IMPROVEMENT OF THE ENHANCED-GPU WITH  
 NVALT IN DIFFERENT QUALITY LOSS

Quality loss	1%	2%	4%	6%	8%	10%
Best CAM size	024-row	512-row	256-row	256-row	128-row	128-row
BlackScholes	1.1x	3.2x	7.3x	10.2x	12.3x	14.7x
FFT	2.7x	6.3x	13.9x	20.5x	23.9x	30.1x
Sobel	2.9x	9.9x	27.7x	38.6x	45.5x	54.0x
inversed2j	1.0x	2.8x	5.7x	8.4x	10.3x	13.3x
Laplacian	3.1x	5.2x	8.4x	12.3x	14.9x	19.3x
Convolution	8.6x	21.2x	32.8x	45.8x	61.4x	79.3x
Binarization	2.9x	5.8x	9.3x	13.5x	17.0x	22.1x
Average	3.2x	7.6x	15.0x	21.3x	26.5x	33.2x

less than 4%. The top-x axis in the figures shows the percentage of time the application runs on CAM for a given THR value. Our evaluation shows that, in small NVALT size, we need to run a larger portion of data on the precise GPU core in order to guarantee the computation accuracy satisfies the required level. By contrast, in large size, the NVALT can provide enough computation accuracy even using very large threshold value. We consider energy-delay product (EDP) to find the optimal CAM size resulting in the best efficiency. Our evaluation shows that 256-row CAM results in the best efficiency of 4.4x energy improvement and 3.4x speedup across these seven test applications (normalized to conventional GPU).

In all applications, 4% quality loss does not show the acceptable quality of service. Therefore, we change the THR value for each application to see the efficiency of the NVALT in different level of accuracy. Table II shows the EDP improvement (normalized to GPU) that our design can achieve for different quality loss from 2% to 10%, as the NVALT size changes. The result shows that our design can tune the level of approximation by partially running the application on NVALT and accurate GPU core. Our result shows that NVALT works much more efficiently for applications which can accept large approximation. For instance, by accepting 6% error, the NVALT can achieve up to 3.7x speedup, while this number decreases to 1.4x to achieve 1% quality loss.

IV. CONCLUSION

This letter suggests a method of hardware approximation that speeds up GPU architecture programs. In order to search for the closest distance row, we develop a (NVALT). NVALT stores high frequency patterns associated with each application, so roughly modeling computing. The stored value that bears the greatest resemblance to the input data is found by our NVALT algorithm. Through the dynamic division of input data between the correct GPU cores and the approximate NVALT, our approach is able to adjust the degree of approximation. NVALT can enhance the GPU, according to our experimental study that runs seven general applications on AMD GPUs.

Compared to the unmodified GPU, energy computation was reduced by 4.5 and performance by 5.7 on average, with less than 10% average relative error.

REFERENCES

- [1] J. Gant and D. Reinsel, "Extracting value from chaos," *IDCview*, vol. 1142, p. 1-12, 2011.
- [2] M. Imani, Y. Kim, A. Rahimi, and T. Rosing, "ACAM: Approximate computing based on adaptive associative memory with online learning," in *Proc. IEEE/ACM ISLPED*, San Francisco, CA, USA, 2016, pp. 162-167.
- [3] A. Yazdanbakhsh, J. Park, H. Sharma, P. Lotfi-Kamran, and H. Esmaeilzadeh, "Neural acceleration for GPU throughput processors," in *Proc. IEEE/ACM Micro*, 2015, pp. 482-493.
- [4] L. Leem, H. Cho, J. Bau, Q. A. Jacobson, and S. Mitra, "ERSA: Error resilient system architecture for probabilistic applications," in *Proc. IEEE/ACM DATE*, Dresden, Germany, 2011, pp. 1560-1565.
- [5] S. Venkataramani, V. K. Chippa, S. T. Chakradhar, K. Roy, and A. Raghunathan, "Quality programmable vector processors for approximate computing," in *Proc. IEEE/ACM Micro*, Davis, CA, USA, 2013, pp. 1-12.
- [6] K. Nepal, Y. Li, R. I. Bahar, and S. Reda, "ABACUS: A technique for automated behavioral synthesis of approximate computing circuits," in *Proc. IEEE/ACM DATE*, Dresden, Germany, 2014, pp. 1-6.
- [7] M. Imani, D. Peroni, and T. Rosing, "CFPU: Configurable floating point multiplier for energy-efficient computing," in *Proc. IEEE/ACM DAC*, Austin, TX, USA, 2017, Art. no. 76.
- [8] M. Imani, A. Rahimi, and T. S. Rosing, "Resistive configurable associative memory for approximate computing," in *Proc. IEEE/ACM*

- DATE, Dresden, Germany, 2016, pp.1327–1332.
- [9] X. Yin, M. Niemier, and X. S. Hu, “Design and benchmarking of ferroelectric FET based TCAM,” in *Proc. IEEE/ACM DATE*, Lausanne, Switzerland, 2017, pp.1444–1449.
- [10] M. Imani, D. Peroni, A. Rahimi, and T. Rosing, “Resistive CAM acceleration for tunable approximate computing,” *IEEE Trans. Emerg. Topics Comput.*, to be published, doi:10.1109/TETC.2016.2642057.
- [11] V. Akhlaghi, A. Rahimi, and R. K. Gupta, “Resistive bloom filters: From approximate membership to approximate computing with bounded errors,” in *Proc. IEEE/ACM DATE*, Dresden, Germany, 2016, pp.1441–1444.
- [12] M. Imani, S. Patil, and T. S. Rosing, “MASC: Ultra-low energy multiple-access single-charge TCAM for approximate computing,” in *Proc. IEEE/ACM DATE*, Dresden, Germany, 2016, pp.373–378.
- [13] M. Imani, S. Patil, and T. Rosing, “Approximate computing using multiple-access single-charge associative memory,” *IEEE Trans. Emerg. Topics Comput.*, to be published, doi:10.1109/TETC.2016.2565262.
- [14] H. Esmailzadeh, A. Sampson, L. Ceze, and D. Burger, “Neural acceleration for general-purpose approximate programs,” in *Proc. IEEE/ACM Micro*, Vancouver, BC, Canada, 2012, pp.449–460.
- [15] R. Ubal, B. Jang, P. Mistry, D. Schaa, and D. Kaeli, “Multi2Sim: A simulation framework for CPU-GPU computing,” in *Proc. ACM PACT*, Minneapolis, MN, USA, 2012, pp.335–344.
- [16] S. Li et al., “McPAT: An integrated power, area, and timing modeling framework for multicore and manycore architectures,” in *Proc. IEEE/ACM Micro*, New York, NY, USA, 2009, pp.469–480.
- [17] S. Kvatinsky, M. Ramadan, E. G. Friedman, and A. Kolodny, “VTEAM: A general model for voltage-controlled memristors,” *IEEE Trans. Circuits Syst. II, Exp. Briefs*, vol.62, no.8, pp.786–790, Aug. 2015.