Stella Nera: Achieving 161 TOp/s/W with Multiplier-free DNN Acceleration based on Approximate Matrix Multiplication

¹Jannis Schönleber, ²Lukas Cavigelli, ²Renzo Andri, ¹Matteo Perotti, ^{1,3}Luca Benini ¹ETH Zürich, Switzerland ²Huawei Technologies, Zürich Research Center, Switzerland ³Università di Bologna, Italy janniss@iis.ee.ethz.ch,{lukas.cavigelli,renzo.andri}@huawei.com,{mperotti,benini}@iis.ee.ethz.ch

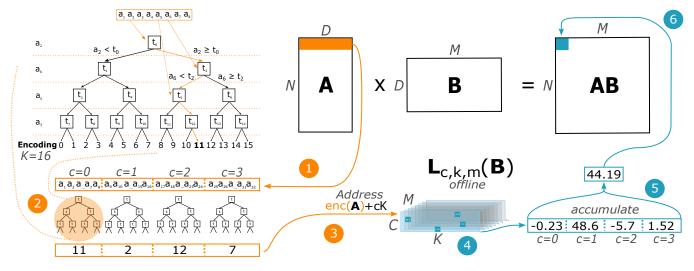


Figure 1: Key steps in Maddness: 1 The input row of the matrix A is divided into C codebooks. 2 Each codebook passes through the encoding function, which maps to the closest prototype. 3 Encoding output and c form the address into the Look-Up Table (LUT) 4. Decoding: the LUT-values from all C codebooks are accumulated 5. Result is written to output AB 6.

ABSTRACT

From classical HPC to deep learning, Matrix Multiplication (MatMul) is at the heart of today's computing. The recent Maddness¹ method approximates MatMul without the need for multiplication by using a hash-based version of product quantization (PQ) indexing into a look-up table (LUT). Stella Nera² is the first Maddness accelerator and it achieves 15x higher area efficiency (GMAC/s/mm²) and a more than 25x higher energy efficiency (TMAC/s/W) than direct MatMul accelerators implemented in the same technology. The hash function is a decision tree, which allows for an efficient hardware implementation as the multiply-accumulate operations are replaced by decision tree passes and LUT lookups. The entire Maddness MatMul can be broken down into parts that allow an effective implementation with small computing units and memories, allowing it to reach extreme efficiency while remaining generically applicable for MatMul tasks. In a commercial 14 nm technology and scaled to 3 nm, we achieve an energy efficiency of 161 TOp/s/W@0.55V with a Top-1 accuracy on CIFAR-10 of more than 92.5% using ResNet9.

²Named after an album of the famous Swiss musicians Patent Ochsner



This work is licensed under a Creative Commons Attribution 4.0 International License.

1 INTRODUCTION

Notwithstanding the end of Dennard's Scaling [12], the latest technology nodes still achieved a $\sim 1.8\times$ better transistor density, a $\sim 1.1\times$ increased maximum frequency, and a $\sim 1.3\times$ power reduction [20] generation over generation. Despite the resulting $\sim 1.45\times$ per-generation energy efficiency improvement on the technology side, the insatiable demand for more computational speed and efficiency needs algorithmic and architectural optimizations to keep up with the fast-paced innovations and ever-growing models for deep learning [32]. As recognized by AMD, the "energy consumption is the primary limiting factor in maintaining the historical rate of performance improvements" [34].

To address this challenge, academia and industry proposed and implemented a multitude of Domain-Specific Architectures (DSA) to accelerate and reduce the power consumption of the most common Machine Learning (ML) workloads. The main optimization target is MatMul, one of the most ubiquitous and critical algorithms that greatly impact overall efficiency and performance. Most of these accelerators exploit quantization to reduce the datapath width, prune the model or exploit sparsity patterns [17, 19, 26, 27] to save on the logic and memory needed for computation and storage, achieving significant improvements in energy efficiency [3, 28, 31]. However, like Moore's Law, quantization is reaching its limits too [4, 29], as Binary Neural Networks (BNNs) encode their weights into bits, the atomic unit of information.

¹Multiply-ADDitioN-lESS

D. Blalock et al. recently presented Maddness [6] as a promising alternative to increase MatMul's efficiency. For the typical ML scenario where one matrix of the MatMul (i.e. the weight matrix) is known ahead of time, Maddness approximates the matrix multiplication by using a hashed version of Product Quantization (PQ) [16] instead of approximating the result by means of a data bit-width reduction. The used hash is a balanced binary regression tree that calculates the following LUT index during the encoding step. Eventually, the accumulation of the looked-up LUT elements yields the result of the MatMul during the following decoding stage. This process replaces all costly Multiply and Accumulate (MACC) operations with inexpensive additions, simplifying the overall MatMul to shifts, adds, and comparisons, paving the way for efficient hardware acceleration.

This work presents the following contributions:

- STELLA NERA, the open-source and fully parameterized implementation of a MADDNESS hardware accelerator. Achieving an energy efficiency of up to 43.1 TOp/s/W@0.55V with a throughput of 2.9 TOp/s@0.55V in a competitive commercial 14nm technology.
- We propose the first formulation of MADDNESS that is differentiable, allowing the method to be used for Deep Neural Networks (DNN) training.
- Fully tested differentiable MADDNESS drop-in Linear and Conv2D layers for PyTorch.
- End-to-end results for ResNet9 [15] on CIFAR-10 reaching an accuracy of 92.6% a difference of 1.2% to the FP32 baseline.

All the results, RTL implementation, verification tests, algorithm implementation, PyTorch layers (including rigorous testing), training scripts, and artefacts are available open-source on https://github.com/joennlae/halutmatmul.

2 RELATED WORK

Our work is influenced by fundamental algorithmic advancements in the encoding function. This function allows us to control the trade-off between accuracy and encoding speed. Matsui et al. provide a comprehensive overview of research directions in the field of product quantization [22, Fig.1]. D. Blalock et al. proposed a fast and efficient encoding function with Bolt [7] by reducing the number of centroids and increasing the number of subvectors compared to previous PQ implementations. However, the similarity function is still norm-based. Maddness improved upon Bolt by replacing the $l_2\text{-norm}$ similarity function (argmin) during the encoding process with a balanced binary regression tree (hash function).

One could also observe steady progress in applying PQ-based methods to DNNs. T. Chen et al. [10] was the first to propose backpropagation by utilizing the hard argmax in the forward path and the differentiable softmax in the backward path to compress embeddings. Ran et al. with PECAN [25] was the first to replace an end-to-end CNN network, also introducing an l_1 -norm-based similarity measure. Tang et al. with LUT-NN [35], and Carter et al. [23] further improved the training methods for end-to-end replaced networks and extended it to BERT, in addition to CNNs, while still using a norm-based encoding function limiting the potential speed-up as one moves Ops from the MatMul to prototype matching in contrast to Maddings where the ceiling is significantly higher.

There is literature on using hashing to accelerate dense layers in neural networks. B. Chen et al. and Spring et al. [9, 11, 33] use locality-sensitive hashing to select which neurons are active in the next layer. These methods significantly differ from Maddness as their goal is to compress and, by that, reduce the number of operations of the network compared to simply approximating the output of the MatMul, allowing Maddness to be used more broadly.

Some of the above-mentioned works provide an optimized CPU implementation using SIMD instructions while Fernandez-Marques et al. [13] propose a custom PQ hardware accelerator (PQA) using an HLS-based implementation on FPGA, although not providing any energy numbers. To the best of our knowledge, we present the first implementation of a MADDNESS accelerator.

3 BACKGROUND

Approximate Matrix Multiplication A general definition for Approximate Matrix Multiplication (AMM) is

$$\|\alpha \cdot dec(enc(\mathbf{A}), l(\mathbf{B})) + \beta - \mathbf{A}\mathbf{B}\|_F < \epsilon(\tau) \|\mathbf{A}\mathbf{B}\|_F,$$
 (1)

where $\mathbf{A} \in \mathbb{R}^{N \times D}$ and $\mathbf{B} \in \mathbb{R}^{D \times M}$ are two matrices. Given $N \gg D \geq M$ and a time budget τ , one has to find the three functions $dec(\cdot, \cdot)$, $enc(\cdot)$ and $l(\cdot)$. For example, DNN meet the requirement for the dimensions N, D, and M. Furthermore, DNN has a weight matrix \mathbf{B} , which is known ahead of time in the inference scenario, a key prerequisite for MADDNESS. N can be seen as the number of images that scale with the training data. $\alpha = 1$ and $\beta = 0$ as normalized data is assumed. To complete the formulation, we assume a training set $\tilde{\mathbf{A}}$, which is drawn from the same distribution as \mathbf{A} . $\|\mathbf{A}\|_F$ refers to the Frobenius norm

$$\|\mathbf{A}\|_{F} = \sqrt{\sum_{i=1}^{N} \sum_{j=1}^{D} \|a_{ij}\|^{2}}.$$
 (2)

Product Quantization PQ is a vector quantization method to approximate inner products or distances for approximate nearest neighbour search (ANN). Conceptually, PQ approximates $ab \approx \tilde{a}b$ for small $\|\tilde{a}-a\|$. \tilde{a} consists of prototypes P that have been learned ahead of time. The prototypes and the known matrix B are multiplied offline to construct the lookup table L using P and B. An encoding function $enc(a) = \tilde{a}$ maps an input vector a to the corresponding prototypes. The resulting prototype corresponds to an offline-learned L entry. For each of the C codebooks (i.e. subspaces), the encoding function maps to one of the K prototypes per subspace. During decoding, the encoded entries of L are summed up to form the result. Vector a is a row of $A \in \mathbb{R}^{N \times D}$ and b is a column of $B \in \mathbb{R}^{D \times M}$.

The essential components of PQ:

• Learning of Prototypes Offline training phase to learn **P**. One splits the input matrix $\tilde{\mathbf{A}}$ into C disjoint subspaces and learns K prototypes via K-means. Use $\tilde{\mathbf{A}}$ as the training set, K prototypes per subspace, and C the number of subspaces. Define $\{\mathbb{I}_c\}_{c=1}^C$ as the indices used in every subspace. The subspaces created by these indices need to be disjoint. Learning prototypes $\mathbf{P} \in \mathbb{R}^{|\mathbb{I}_c| \times K \times C}$ means minimizing

$$\sum_{n=1}^{N} \sum_{c=1}^{C} \sum_{i \in \mathbb{I}_{c}} (\tilde{\mathbf{A}}_{ni} - \mathbf{P}_{i,\mathbf{z},c})^{2}.$$
 (3)

One can achieve that by running K-means in every subspace \mathbb{I}_c to learn the cluster mapping $\mathbf{z}_{n,c,i} \in \mathbb{R}^{\mathbb{N}} = \mathbf{z}$ and write the resulting

centroids into $\mathbf{P}_c \in \mathbb{R}^{|\mathbb{I}_c| \times K}$ where $|\mathbb{I}_c| = \frac{D}{C}$.

• Encoding Function enc(A,P) 2 A function that calculates the closest prototype based on a similarity measurement. This is done for every subspace C, resulting in C encoding indices per input row. The default similarity measurement l_2 -norm results in the encoding function

$$enc_c(\mathbf{a}) = \underset{k}{\operatorname{argmin}} \sum_{i \in \mathbb{T}_c} (\mathbf{a}_i - \mathbf{P}_{i,k,c})^2. \tag{4}$$

• LUT $L \in \mathbb{R}^{M \times C \times K}$ 4 is computed offline using the determined prototypes $\mathbf{P} \in \mathbb{R}^{\frac{D}{C} \times K \times C}$ and weight matrix $\mathbf{B} \in \mathbb{R}^{D \times M}$. Let $\tilde{\mathbf{P}} \in \mathbb{R}^{D \times K \times C}$ be the expanded version of \mathbf{P} where one adds zero for all indices in the first dimension that are not in the corresponding subspace $i \notin \mathbb{I}_c$. After reordering the dimensions $\tilde{\mathbf{P}} \in \mathbb{R}^{D \times K \times C} \to \tilde{\mathbf{P}} \in \mathbb{R}^{C \times K \times D}$.

$$lut(\mathbf{B}^{T}, \tilde{\mathbf{P}}) = \mathbf{L}_{m,c,k} = \sum_{i=1}^{D} \mathbf{B}_{m,i} \mathbf{P}_{c,k,i}.$$
 (5)

 Accumulation dec(enc(a, P), L) sums over C lookup table entries with the determined indices from the encoding step,

$$(\mathbf{A}\mathbf{B})_{n,m} \approx \sum_{c=1}^{C} \mathbf{L}_{m,c,k}, k = enc_c(\mathbf{a}_n).$$
 (6)

3.1 MADDNESS

MADDNESS replaces two components of PQ: the encoding function $enc_c(\mathbf{a})$ and the values stored in the LUT L (\mathbf{a} Fig. 1). A balanced binary tree (\mathbf{a} Fig. 1) replaces the hash function and the results of the vector-vector multiplication between the prototypes \mathbf{P} and the matrix \mathbf{B} are stored inside \mathbf{L} . Using a hash function instead of the l_2 -norm $\|a-b\|_2$ results in a significant speedup with a slight accuracy loss. This allows MADDNESS to be used on large input matrices \mathbf{A} as every input row has to be encoded.

4 DIFFERENTIABLE MADDNESS

A significant limitation of the hash function is its non-differentiable nature, making it unusable in gradient-based optimizations. This reformulation builds on the previous work by Zhang [37] for decision trees.

Let S_c \in $\{0,1\}^{15\times4}$ be the selection matrix and H_c \in $\{-1,0,1\}^{16\times15}$ the tree description matrix.

$$S = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{bmatrix} \quad H = \begin{bmatrix} -1 & -1 & 0 \\ -1 & 1 & 0 \\ 1 & 0 & -1 \\ 1 & 0 & 1 \end{bmatrix} \quad \begin{bmatrix} t_0 \\ t_1 \\ t_2 \\ / / / \\ 0 & 1 & 2 & 3 \end{bmatrix}$$

Figure 2: Example: S/H for the small two-level decision tree.

To parallelize over multiple codebooks, let:

$$\mathbf{S} = \bigoplus_{i=1}^{C} \mathbf{S}_i = \operatorname{diag}(\mathbf{S}_1, \mathbf{S}_2, ..., \mathbf{S}_C), \quad \mathbf{H} = \bigoplus_{i=1}^{C} \mathbf{H}_i = \operatorname{diag}(\mathbf{H}_1, \mathbf{H}_2, ..., \mathbf{H}_C) \quad (7)$$

Let $\theta \in \mathbb{R}^{C \times 15}$ be the decision threshold vector. Then, we define the encoding matrix

$$E = \operatorname{argmax}(H\sigma(SX - \theta)), \tag{8}$$

where $\sigma=$ sign function. Equation 8 is for the forward pass. To calculate the gradient, we formulate

$$\mathbf{E}_{soft} = \operatorname{softmax}(\mathbf{H}\sigma(\mathbf{SX} - \theta)) \tag{9}$$

using σ = tanh. This allows us to pass the gradient back through the decision tree by employing the technique of straight-through estimator (STE) [5].

Using $E \in \{0,1\}^{N \times C \times K}$ and the LUT L (from (5)), we first sum over K, then accumulate over C:

$$\mathbf{I}_{n,m,c} = \sum_{i=1}^{K} \mathbf{E}_{n,c,i} \mathbf{L}_{m,c,i}, \qquad (\mathbf{AB})_{n,m} = \sum_{i=1}^{C} \mathbf{I}_{n,m,c}.$$
 (10)

Using this approach to backpropagate the gradient, we provide a PyTorch implementation for Linear and Conv2D layers. This approach can also be used to increase accuracy after the offline MADDNESS initialization or to start from a random initialization.

Before applying Maddness, Conv2D layers use the im2col method [8] to convert the Conv2D to a MatMul using input $\mathbf{X} \in \mathbb{R}^{N \times C_i \times H_i \times W_i}$, weights $\mathbf{W} \in \mathbb{R}^{C_o \times C_i \times h_k \times w_k}$, the im2col transformed input $\tilde{\mathbf{X}} \in \mathbb{R}^{N \cdot H_i \cdot W_i \times C_i \cdot h_k \cdot w_k}$ is defined as

$$\tilde{\mathbf{X}}(b\cdot N+o_x\cdot W_i\cdot o_y,i_c\cdot h_k\cdot w_k+k_x\cdot w_k+k_y)=\mathbf{X}(b,i_c,o_x+k_x,o_y+k_x)$$
 for padding and stride equals 1 with $b=1,...,N,i_c=1,...,C_i,o_x$ and $o_y=1,...,H_i$ and W_i,k_x and $k_y=1,...,h_k$ and $w_k.\tilde{\mathbf{W}}\in\mathbb{R}^{C_i\cdot h_k\cdot w_k\times C_o}$ are the transformed weights. Resulting in the Conv2D being represented as $\mathbf{O}=\tilde{\mathbf{X}}\cdot\tilde{\mathbf{W}}$ a MatMul with an inner dimension D of $C_i\cdot h_k\cdot w_k$ allowing us to learn one codebook per input channel when using a codebook width of $h_k\cdot w_k$. Whereas for kn2col $(H_i\cdot W_i\times C_i)\times (C_i\times k_x\cdot k_y\cdot C_o)$ with an inner dimension of C_i using the same codebook-width, keeping the LUT size equal, results in lower accuracy due to the resulting prototypes being across channels.

To learn the LUT in INT8, we employ another STE where the INT8 LUT is used during the forward pass and, in the backward pass the floating-point version of the LUT. After each backward pass, the INT8 LUT is requantized.

5 ARCHITECTURE

We will describe the architecture from the ground up, beginning with encoding and decoding units and gradually building up to a complete system capable of handling a ResNet-like architecture.

Encoding Unit To encode an input vector, the encoding unit moves through the binary balanced tree. It takes four cycles to generate a valid encoding, as it needs to go through all four levels of the tree (K = 16). Each Stella Nera accelerator implementation has four encoders that build a unit. These encoders operate with a one-tree-level/cycle offset to produce one valid encoding per cycle together. The depth of the tree (i.e. the number of prototypes \sqrt{K}) is a parameter of the Stella Nera architecture. K = 16 results in the best accuracy vs. performance trade-off [35].

Decoding Unit The LUT used by Stella Nera is built with standard cell-based memory (SCM), and it processes one encoded value or address into the LUT per cycle. The value obtained from the LUT is then added to the current accumulation. After C accumulations (6),

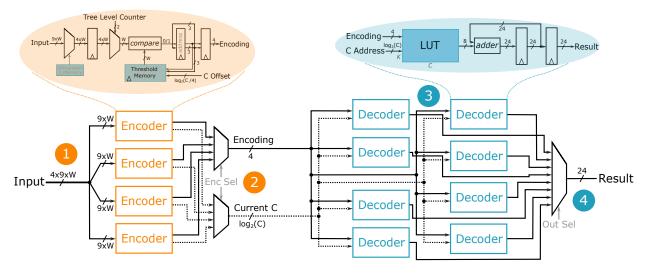


Figure 3: Architecture of STELLA NERA: 1 The input gets split up into the four encoders. The encoders are offset by one cycle and need four cycles to traverse a tree. This leads to one valid output per cycle 2. The encoded value and the current C are the LUT address 3. After C cycles, all decoders have a valid result. A multiplexer selects the results in order 4.

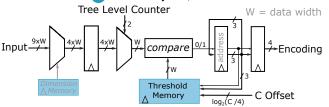


Figure 4: Schematic of the Encoding Unit

the output value is the accumulated value. Two variants were implemented, namely, INT8 LUT and INT24 accumulation and FP16/FP32. As the LUT is already approximate, the use of INT8 results in a negligible accuracy drop, as also found by Tang et al. [35]. Therefore, we will only report the results of the INT8/INT24 version.

System & Scaling When we combine encoder and decoder, as illustrated in Fig. 3, we get the Stella Nera accelerator. The schematic shows four encoders and eight decoders, with four encoders producing one output per cycle.

Stella Nera has two scaling parameters. The first is the total number of decoders called N_{dec} , while the second is the number of codebooks, C_{dec} , per decoder. It is important to note that the subparameter W_{dec} encompasses the number of output values per cycle, which is necessary when N_{dec} is greater than C_{dec} . In such a case, serving all outputs in the C_{dec} cycles required to calculate the next value would not be possible.

The most effective approach is to combine multiple smaller units based on the given problem size. For example, if we have four Stella Nera units with $N_{dec}=64$, $C_{dec}=16$ and $W_{dec}=8$, we can tile in all three tiling directions: N,D,M, providing maximum flexibility. However, an additional adder is required if we tile in the D (i.e., C) direction.

6 ALGORITHMIC RESULTS

We evaluated our differentiable implementation of MADDNESS using the ResNet9 architecture with the channel numbers [64, 128, 128,

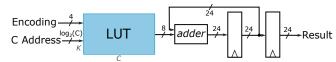


Figure 5: Schematic of the Decoding Unit

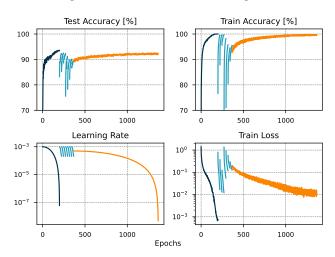


Figure 6: Key metrics for the ResNet9 training colored with the three different stages pre-training, layer-by-layer and fine-tuning.

256,256,256,256], which is a commonly used benchmark architecture for edge scenarios in deep learning. The dataset used for training and evaluation was CIFAR-10. A base model was trained for $N\!=\!200$ epochs with a batch size of 128. For fine-tuning with MADDNESS, we replaced the layers using a codebook width of $CW\!=\!9$ and a number of prototypes $K\!=\!16$. This resulted in a LUT table size twice the size of the weights and one unrolled $3\!\times\!3$ kernel per codebook. The

	ISSCC	ISSCC	ISSCC	Nature Elec.	VLSI Symp.	DATE	This work	This work
	2022 [24]	<i>2020</i> [18]	<i>2021</i> [1]	2023 [14]	2018 [2]	<i>2022</i> [36]	4 units	$scaled^5$
Technology [nm]	4	7	7	14	14	22	14	3
Power Supply [V]	0.55	0.575	0.55	0.85	0.55	0.65	0.55	0.55
Precision	INT8 ¹	INT8 ¹	INT4 ¹	12b ADC	INT8	FP16	INT8 (LUT)	INT8 (LUT)
ResNet9 Acc. [%]	93.6 ⁷	93.6 ⁷	91.7 ⁷	92.8 ⁶	93.6 ⁷	93.7	92.6	92.6
Peak Energy Eff. [TOp/s/W]	11.6 ²	6.8	16.5	1.7	5.2	0.7	43.14	161 ⁴
Peak Area Eff. [TOp/s/mm ²]	3.45^{2}	0.4^{3}	3.3	0.34	1.7 ³	0.06	5.1 ⁴	163 ⁴
Peak Perf. [TOp/s]	19.7	1.2 ³	64	1.7	0.04^{3}	0.03	2.9^{4}	4.1 ⁴
Area [mm ²]	4.74	3.04	19.6	11.18	0.03	0.5	0.57	0.025
Frequency [MHz]	332	290	1'000	-	400	476	624	886
Power [mW]	381 ²	174	3'880 ³	-	8	43.5	60.9	23.0

¹ These chips support additional precisions. ³ Inferred numbers. ⁴ using a codebook width CW = 9. ⁵ scaled using [30] and [20]

Table 1: Comparison of state-of-the-art. 1 MAC = 2 Ops.

replaced layers were initialized using the Maddness algorithm [6, Algorithm 2]. For each layer, the network was trained for N=25 epochs with a learning rate $\eta=10^{-3}$ and half of that learning rate for the thresholds using a cosine annealing scheduler with a $T_{max}=25$ and a $\eta_{min}=2\cdot 10^{-4}$.

We followed the common practice of calculating the first and last layer in FP16. In our case, these are less than 1% of compute operations. Once all layers were replaced, the network was trained for N=1000 epochs with a learning rate $\eta=5\cdot 10^{-4}$ using a cosine annealing scheduler with $T_{max}=1000$. We employed commonly used tactics such as random cropping and random horizontal flipping. The pre-training, layer-by-layer and fine-tuning stages are visualized in Fig. 6. We achieved an end-to-end accuracy of 92.64% (see Table 1).

7 HARDWARE EVALUATION

Physical Implementation Our accelerator was implemented using industry-standard EDA tools and a 14nm technology. All the numbers presented are taken from an SDF backannotated activity-based post-layout simulation (TT, 0.55V, 25°C) for energy or extracted from the post-layout design running a MatMul using randomly generated inputs. The STELLA NERA accelerator runs at 624 MHz (using a maximum of five levels of logic), accounting for hold and setup timing in the respective worst-case conditions. The floorplan of the STELLA NERA accelerator, which uses N_{dec} = 64, W_{dec} = 8, C_{dec} = 16, is shown in Fig. 7. As expected, a significant portion of the area is used by LUTs, shown with the numbers. One can also spot the four encoding units in and INT8/INT24 adders in .

Subunit Analysis Table 2 gives an overview of each subunit of the Stella Nera accelerator during the Maddness execution. The energy needed per LUT lookup is 0.23 pJ, significantly higher than the negligible amount of energy consumed by the INT8/INT24 adder. This means that inexpensive tiling along the D (i.e., C) direction of the MatMul can be achieved by adding an additional adder when

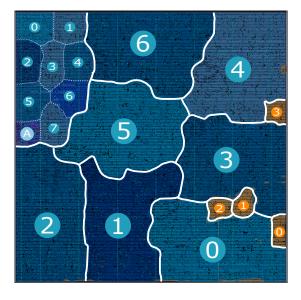


Figure 7: Physical implementation of a STELLA NERA accelerator with unit N_{dec} = 64, W_{dec} = 8, C_{dec} = 16 using a commercial 14nm technology. The core area for this implementation is 350 μ m \times 350 μ m at a utilization of 85% (860 kGE).

using multiple Stella Nera units. One decoding operation involves CW MACs or $CW \cdot 2$ Ops, resulting in an Energy/Op of around 30 fJ for CW = 9. The encoding energy is divided by all decoders P_{dec} as they use the same encoding.

System Results For our proposed system, we combine four Stella Nera units ($P_{dec} = 64, W_{dec} = 64, C_{dec} = 16$) together. Expect for twice the size of the "weights" (for CW = 9), i.e., LUT, the bandwidth needs compared to the standard MatMul stay the same, and no weight buffer is needed as the weights are stored inside the

² End-to-end energy numbers running MobileNetEdgeTPU. Their peak energy efficiency is higher.

⁶ They use slightly different channel numbers to fit their specific architecture.

⁷ Trained via quantization-aware training (QAT) using PyTorch, the code is available online. ⁸ Area for 8 cores.

Unit	Area	GE		Energy			
				Dynamic	Ste	$per\ Op$	
	$[\mu m^2]$	$[\mathrm{kGE}]$	[mW]	Switching	Leaking	Internal	[pJ]
Encoder	790	6.5	0.20	68.3 %	3.4 %	28.2 %	0.33
Encoder 4x	3'190	26.4	0.91	69.8 %	3.1 %	27.1 %	
Decoder	1'570	13.0	0.16	57.3 %	6.0 %	36.7 %	0.26
Decoder 8x	12'650	104.6	1.48	60.6 %	5.3 %	34.1 %	2.37
INT8/INT24 Adder	35	0.3	0.02	63.4 %	1.2 %	35.4 %	0.03

Table 2: Subunit area, power, and energy per operation during the execution of MADDNESS in 14nm, TT, LVT, 0.55V.

accelerator. We compare with the state-of-the-art in Table 1. We also provide scaled numbers for the newest technology we compare, using DeepScaleTool [30] from 14nm to 7nm and foundry-published scale numbers from 7nm to 3nm [20]. As our current implementation uses FP16 as input and will output FP16 in the end, we need an INT8 to FP16 conversion, a scaling plus a shift resulting in an FMA operation. To account for that, we add two FP16 FMA units [21] per unit, each capable of two conversions per cycle. This adds 13.7 pJ/cycle or 8.5 mW (23.3%) to the four units STELLA NERA system. Resulting in 984 inf/sec @ 23.7uJ per inference (in 14 nm) for ResNet9, where of 9.2uJ are spent in the first nonaccelerated FP16 layer.

8 FUTURE WORK

While our proposed differentiable Maddness implementation has demonstrated promising results, the model size of ResNet9 is around eight times larger in parameter size compared to cutting-edge models. Some encouraging directions are model architectures that are more suitable to Maddness, eliminating the need for INT8 to FP16 conversion by quantizing the input and thresholds, temperature parameters for the STEs to reach higher accuracy, and connecting the encoder and adders directly to an existing SRAM in a system.

9 CONCLUSION

We present the first differentiable implementation of MADDNESS unlocking DNN training. We provide the necessary training scripts and artefacts for a 92.5%+ accuracy ResNet9 model using MADDNESS. We introduce STELLA NERA, a fully parametrizable, open-source MADDNESS accelerator implemented in a 14nm technology. It reaches an energy efficiency of 43.1 TOp/s/W@0.55V in 14 nm or 161 TOp/s/W@0.55V scaled to 3 nm with a peak performance of 2.9 TOp/s and 4.1 TOp/s respectively, resulting in 15x higher area efficiency and 25x higher energy efficiency than the state-of-the-art. Code, RTL and scripts are available on https://github.com/joennlae/halutmatmul.

REFERENCES

- Ankur Agrawal et al. 2021. 9.1 A 7nm 4-core AI chip with 25.6 TFLOPS hybrid FP8 training, 102.4 TOPS INT4 inference and workload-aware throttling. In *Proc. IEEE ISSCC*, Vol. 64. 144–146.
- [2] Mark Anders et al. 2018. 2.9 tops/w reconfigurable dense/sparse matrix-multiply accelerator with unified int8/inti6/fp16 datapath in 14nm tri-gate cmos. In 2018 IEEE Symposium on VLSI Circuits. IEEE, 39–40.
- [3] Renzo Andri et al. 2018. Hyperdrive: A Systolically Scalable Binary-Weight CNN Inference Engine for mW IoT End-Nodes. In Proc. IEEE ISVLSI. 509–515.
- [4] Renzo Andri et al. 2021. ChewBaccaNN: A Flexible 223 TOPS/W BNN Accelerator. In Proc. IEEE ISCAS. https://doi.org/10.1109/ISCAS51556.2021.9401214
- [5] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. 2013. Estimating or Propagating Gradients Through Stochastic Neurons for Conditional Computation. arXiv:1308.3432 [cs.LG]
- [6] Davis Blalock and John Guttag. 2021. Multiplying matrices without multiplying. In International Conference on Machine Learning. PMLR, 992–1004.

- [7] Davis W Blalock and John V Guttag. 2017. Bolt: Accelerated data mining with fast vector compression. In Proc. ACM KDD. 727–735.
- [8] Lukas Cavigelli, Michele Magno, and Luca Benini. 2015. Accelerating Real-Time Embedded Scene Labeling with Convolutional Networks. In Proc. ACM/IEEE Design Automation Conference.
- [9] Beidi Chen, Tharun Medini, and Anshumali Shrivastava. 2019. SLIDE: In Defense of Smart Algorithms over Hardware Acceleration for Large-Scale Deep Learning Systems. arXiv:1903.03129 (2019).
- [10] Ting Chen, Lala Li, and Yizhou Sun. 2020. Differentiable Product Quantization for End-to-End Embedding Compression. In Proc. ICML. 1617–1626.
- [11] Wenlin Chen et al. 2015. Compressing Neural Networks with the Hashing Trick. In Proceedings of the 32nd International Conference on Machine Learning. PMLR.
- [12] R.H. Dennard et al. 1974. Design of ion-implanted MOSFET's with very small physical dimensions. IEEE Journal of Solid-State Circuits 9, 5 (1974), 256–268.
- [13] Javier Fernandez-Marques et al. 2023. Are We There Yet? Product Quantization and its Hardware Acceleration. arXiv:2305.18334 [cs.AR]
- [14] Manuel Le Gallo et al. 2023. A 64-core mixed-signal in-memory compute chip based on phase-change memory for deep neural network inference. *Nature Electronics* 6, 9 (Aug. 2023), 680–693.
- [15] Kaiming He et al. 2016. Deep residual learning for image recognition. In Proceedings of the IEEE conference on computer vision and pattern recognition. 770–778.
- [16] Herve Jégou, Matthijs Douze, and Cordelia Schmid. 2011. Product Quantization for Nearest Neighbor Search. IEEE Trans. on Patt. Analysis and Machine Int. (2011).
- [17] Behnam Khaleghi et al. 2022. PatterNet: Explore and Exploit Filter Patterns for Efficient Deep Neural Networks. In Proc. ACM/IEEE DAC.
- [18] Chien-Hung Lin et al. 2020. 7.1 A 3.4-to-13.3 TOPS/W 3.6 TOPS dual-core deep-learning accelerator for versatile AI applications in 7nm 5G smartphone SoC. In *Proc. IEEE ISSCC*. 134–136.
- [19] Xiaolong Lin et al. 2023. AdaS: A Fast and Energy-Efficient CNN Accelerator Exploiting Bit-Sparsity. In Proc. ACM/IEEE DAC.
- [20] Mark Liu. 2021. 1.1 Unleashing the Future of Innovation. In 2021 IEEE International Solid-State Circuits Conference (ISSCC), Vol. 64. 9–16.
- [21] Stefan Mach et al. 2020. FPnew: An open-source multiformat floating-point unit architecture for energy-proportional transprecision computing. IEEE Transactions on Very Large Scale Integration (VLSI) Systems 29, 4 (2020), 774–787.
- [22] Yusuke Matsui et al. 2018. A survey of product quantization. ITE Trans. on Media Technology and Applications 6, 1 (2018), 2–10.
- [23] Calvin McCarter and Nicholas Dronen. 2022. Look-ups are not (yet) all you need for deep learning inference. arXiv preprint arXiv:2207.05808 (2022).
- [24] Jun-Seok Park et al. 2022. A Multi-Mode 8K-MAC HW-Utilization-Aware Neural Processing Unit with a Unified Multi-Precision Datapath in 4nm Flagship Mobile SoC. In Proc. IEEE ISSCC, Vol. 65, 246–248.
- [25] Jie Ran et al. 2023. PECAN: A Product-Quantized Content Addressable Memory Network. In Proc. IEEE DATE.
- [26] Enrico Reggiani, Renzo Andri, and Lukas Cavigelli. 2023. Flex-SFU: Accelerating DNN Activation Functions by Non-Uniform Piecewise Approximation. In Proc. ACM/IEEE Design Automation Conference (DAC).
- [27] Ao Ren et al. 2023. HBP: Hierarchically Balanced Pruning and Accelerator Co-Design for Efficient DNN Inference. In Proc. ACM/IEEE DAC.
- [28] Albert Reuther et al. 2022. AI and ML Accelerator Survey and Trends. In 2022 IEEE High Performance Extreme Computing Conference (HPEC). 1–10.
- [29] Manuele Rusci, Lukas Cavigelli, and Luca Benini. 2018. Design Automation for Binarized Neural Networks: A Quantum Leap Opportunity?. In 2018 IEEE International Symposium on Circuits and Systems (ISCAS).
- [30] Satyabrata Sarangi and Bevan Baas. 2021. DeepScaleTool: A Tool for the Accurate Estimation of Technology Scaling in the Deep-Submicron Era. In Proc. IEEE ISCAS.
- [31] Moritz Scherer et al. 2022. CUTIE: Beyond PetaOp/s/W Ternary DNN Inference Acceleration With Better-Than-Binary Energy Efficiency. IEEE Trans. on Computer-Aided Design of Integrated Circuits and Systems 41, 4 (2022), 1020–1033.
- [32] Jaime Sevilla et al. 2022. Compute Trends Across Three Eras of Machine Learning. In Proc. International Joint Conference on Neural Networks (IJCNN).
- [33] Ryan Spring and Anshumali Shrivastava. 2017. Scalable and Sustainable Deep Learning via Randomized Hashing. In Proceedings of the 23rd ACM SIGKDD.
- [34] Lisa Su and Sam Naffziger. 2023. 1.1 Innovation For the Next Decade of Compute Efficiency. In *Proc. IEEE ISSCC*. 8–12.
- [35] Xiaohu Tang et al. 2023. LUT-NN: Empower Efficient Neural Network Inference with Centroid Learning and Table Lookup. In Proc. ACM MobiCom.
- [36] Yvan Tortorella et al. 2022. RedMulE: A Compact FP16 Matrix-Multiplication Accelerator for Adaptive Deep Learning on RISC-V-Based Ultra-Low-Power SoCs. In Proc. IEEE DATE.
- [37] Jinxiong Zhang. 2022. Yet Another Representation of Binary Decision Trees: A Mathematical Demonstration. arXiv:2101.07077 [cs.LG]