POLITECNICO DI TORINO Repository ISTITUZIONALE

MARLIN: A Co-Design Methodology for Approximate ReconfigurabLe Inference of Neural Networks at the Edge

Original

MARLIN: A Co-Design Methodology for Approximate ReconfigurabLe Inference of Neural Networks at the Edge / Guella, Flavia; Valpreda, Emanuele; Caon, Michele; Masera, Guido; Martina, Maurizio. - In: IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS. I, REGULAR PAPERS. - ISSN 1549-8328. - ELETTRONICO. - (2024). [10.1109/tcsi.2024.3365952]

Availability: This version is available at: 11583/2986451 since: 2024-02-29T11:26:20Z

Publisher: IEEE

Published DOI:10.1109/tcsi.2024.3365952

Terms of use:

This article is made available under terms and conditions as specified in the corresponding bibliographic description in the repository

Publisher copyright

(Article begins on next page)

MARLIN: a Co-design <u>Methodology</u> for <u>Approximate ReconfigurabLe Inference of Neural</u> Networks at the Edge

Flavia Guella[®] Student member, IEEE, Emanuele Valpreda[®] Student member, IEEE, Michele Caon[®] Student member, IEEE, Guido Masera[®] Senior member, IEEE, Maurizio Martina[®] Senior member, IEEE

Abstract—The optimization of neural networks (NNs) is nec-1 essary to enable their deployment on energy-constrained devices. 2 State-of-the-art methods leverage approximate multipliers to 3 execute NNs reducing the inference energy without heavily 4 affecting the accuracy. However, previous works usually require 5 a specialized hardware accelerator and are limited to fixed multi-6 pliers or reconfigurable ones with few approximation levels. This 8 paper introduces MARLIN, a framework to deploy layerwise approximate NNs on PULP, a microcontroller with a RISC-V 9 core. A multiplier architecture, with runtime selection of 256 10 approximation levels, is developed and integrated into the PULP 11 cluster cores, enabling runtime configuration through control 12 status register (CSR) instructions embedded within the code. 13 The PULP toolchain is adapted to incorporate the approximation 14 level selection within the instruction flow seamlessly. MARLIN 15 leverages the genetic algorithm NSGA-II to search for the 16 17 best configurations among thousands of approximate NNs. The framework is validated by simulating an approximate NN trained 18 with the MNIST dataset on PULP. Moreover, MARLIN is used 19 to optimize and approximate six ResNet models trained with the 20 CIFAR-10 dataset. In particular, for ResNet-56, the most complex 21 NN used in the experiments, the multiplication energy is reduced 22 by 23.9% while retaining 99% of the accuracy of the exact model. 23 24

Index Terms—Approximate computing, neural networks,
 RISC-V, hardware acceleration, reconfigurable computing.

I. INTRODUCTION

27

²⁸ IN recent years, there has been an increasing necessity
 ²⁹ to include neural networks (NNs) in embedded systems
 ³⁰ to deliver more advanced functionalities. Nonetheless, NNs
 ³¹ superior task accuracy comes at the cost of high computational

Flavia Guella, Emanuele Valpreda, Michele Caon, Guido Masera, and Maurizio Martina are with Department of Electronics and Telecommunications, Politecnico di Torino, Torino, 10129, Italia (email: flavia.guella@polito.it, emanuele.valpreda@polito.it).

complexity and memory requisites, thereby presenting chal-32 lenges in deploying them on energy-constrained devices, such 33 as microcontrollers (MCUs) [1], [2]. Over the past decade, 34 dedicated hardware accelerators [3]-[9] have been developed 35 to optimize energy efficiency and throughput during inference 36 by reducing the data movement and the cost of arithmetic 37 operations. The latter usually accounts for around a fourth 38 of the total inference energy, as the convolutional and fully 39 connected layers of NNs involve millions of multiplications 40 and additions [10]. Nowadays, quantization-aware training 41 can reduce the bitwidth of NN models to 8 bits or below 42 with little or no loss in accuracy [3], [11]–[13]. It lowers 43 the amount of memory traffic and the computational cost, 44 allowing the deployment of NNs on low-power devices with no 45 support for floating-point arithmetic. Layer-wise quantization 46 leverages the different degrees of robustness and tolerance 47 to error introduction of NN layers. Therefore, a runtime 48 reconfigurable multiplier supporting operands with different 49 bitwidths is necessary to leverage mixed-precision and layer-50 wise quantization [3], [4], [11], [12]. Similarly to quantization, 51 approximation is another co-design technique mainly aimed 52 at reducing the inference energy [14]. The basic idea is to 53 reduce computational complexity and cost using operators that 54 produce inexact results. However, as already pointed out, there 55 is high variability in the sensitivity of single layers inside the 56 same model and among different models. Therefore, designing 57 a multiplier with different approximation levels is fundamental 58 to ensure flexibility and adaptability. Several strategies have 59 been explored in the literature to design hardware supporting 60 layer-wise approximate NNs [15]-[17], leveraging retraining 61 or parameters fine-tuning to reduce the accuracy degradation. 62 However, previous works rely on specialized accelerators and 63 support very few approximation levels [17], [18], limiting 64 the possibility of finding the optimal error-accuracy trade-65 off; other works use several non-reconfigurable multipliers 66 instances in the accelerator's systolic array [15], [19], [20]. 67 The latter approach increases the area overhead and prevents 68 the same hardware from executing NNs with different pa-69 rameters or NNs with different model architectures with the 70 same energy efficiency or accuracy. Since Internet of Things 71 (IoT) devices have limited area and power budgets [1], [2] and 72 must be able to adapt to different workloads and performance 73 targets, it is necessary to adopt a layer-wise approximation 74 strategy that relies on a single multiplier offering several ac-75 curacy levels. Moreover, the architecture using this multiplier 76

This article has been accepted for publication in IEEE Transactions on Circuits and Systems--I: Regular Papers. This is the author's version which has not been fully edited and content may change prior to final publication. Citation information: DOI 10.1109/TCSI.2024.3365952

JOURNAL OF LATEX CLASS FILES, VOL. 18, NO. 9, SEPTEMBER 2020

2

132

133

⁷⁷ should be programmable and reconfigurable to support the
⁷⁸ deployment of different NNs, reusing the same hardware as
⁷⁹ efficiently as possible without requiring a redesign or recall of
⁸⁰ the IoT device. Such flexibility allows choosing which NN is
⁸¹ suited for a particular scenario, prioritizing the battery life of
⁸² the edge device or the accuracy of the predictions.

In this work, we present MARLIN, an automated layer-wise 83 approximation co-design methodology that enables searching 84 for the optimal energy-accuracy trade-off and deploying ap-85 proximate NNs on a hardware accelerator. We leverage a 86 runtime-reconfigurable parallel tree multiplier featuring 256 87 error levels, assigning low-power inaccurate multiplier con-88 figurations to layers with high error resilience and more 89 accurate, albeit less efficient, setup to layers with low error 90 tolerance. Moreover, we search for different network-level 91 approximate configurations, i.e., NNs with the same model 92 architecture but a different layer-wise approximation and, 93 consequently, energy-accuracy trade-off, supporting different 94 workload priorities. In order to rapidly deploy an approximate 95 NN and prove the portability of this method to a known open-96 source hardware, we selected the PULP MCU [21] as the 97 target IoT platform with RI5CY [22], a RISC-V core. The 98 search of the approximate NNs configuration is done with a 99 non-dominated sorting genetic algorithm (NSGA-II) [23]. This 100 procedure is automated and executed offline, before the NN 101 is mapped to any specific hardware. The configuration of the 102 arithmetic units, using the search results, is done online. The 103 contributions of this work are summarized as follows: 104

A layer-wise approximation strategy that reduces the en-105 ergy of arithmetic operations while retaining the original 106 task accuracy, applicable to any NN topology with convo-107 lutional and fully connected layers, with no modification 108 to the model architecture. During the optimization pro-109 cess, relying on NSGA-II, approximate NNs are evaluated 110 by their error resilience and energy, assigning to each 111 layer a different multiplier configuration. 112

- A runtime reconfigurable signed multiplier supporting 256 approximation levels.
- The post-synthesis simulation of the proposed methodology on a RI5CY core, adapted to include the proposed approximate multiplier, supporting the runtime selection of the accuracy. The entire deployment process of the approximate NN is automated and requires no specific knowledge of the optimization method.
- The code used in this work, the experimental data, and the instructions to replicate the results presented in this paper are available at https://github.com/vlsi-lab/MARLIN

This article is organized as follows: Section II introduces a 124 background on NN accelerators and relevant works on ap-125 proximate inference, discussing the differences with MARLIN, 126 Section III details the co-design framework, the search strat-127 egy, the multiplier design, and the modified PULP-toolchain, 128 Section IV presents the results, a comparison with the state-129 of-the-art, and discusses the trade-offs of reconfigurable ap-130 proximate computing, and Section V summarizes the paper. 131

II. BACKGROUND AND RELATED WORKS

A. Hardware Accelerators and Mapping

NN accelerators are typically based on a hardware archi-134 tecture that comprises a memory hierarchy and an array of 135 interconnected processing engines (PEs) [4], [7]-[10], similar 136 to the one depicted in Figure 1. The memory hierarchy usually 137 comprises the system's main memory (off-chip), global buffers 138 (on-chip), and the registers within the processing engines. 139 The off-chip and on-chip memories are connected through 140 the system's bus, whereas the PEs communicate through a 141 network-on-chip. The execution of a NN is scheduled with a



Fig. 1. Generic high-level model of a hardware accelerator for NN inference.

142 dedicated mapper that generates the instructions and partitions 143 the resources, minimizing the energy and latency associated 144 with the data movement. The complete cost of a multiplication 145 considered by a mapper includes the energy and latency 146 required to read all the operands, move them through the 147 memory hierarchy, compute and store the result [9], [10], [24], 148 [25]. The mappers used in [7]–[9] can be compared to the com-149 pilers used to generate machine code for processors such as 150 RI5CY [22], which purpose is to optimize the performance and 151 resource usage. The approximation methodology described in 152 Section III is orthogonal to the mapping process as it only 153 modifies the energy associated with the arithmetic operations 154 and not the data movement. In [24] a tool that predicts the 155 energy of approximate NN with a mapper based on [10] is 156 presented validating the assumption that approximate comput-157 ing only modifies the arithmetic energy without affecting the 158 data movement. Therefore, MARLIN does not influence the 159 mapping and could be easily integrated in [7]-[9], similarly to 160 what has been done for DORY in Section III-E, as discussed in 161 Section IV-D. Alternative hardware architectures for energy ef-162 ficient NN inference at the edge are MCUs. ISA extensions and 163 dedicated dot-product units inserted in the RISC-V pipeline are 164 used in [5], [6], whereas in [26] a low-power neural processing 165 unit is connected to the MCU through the bus. In this work, 166 we selected PULP [5] as the target platform, a MCU-based 167 low-power computing platform with a host CPU relying on 168 the RI5CY or zero-riscy cores and equipped with a multi-169 core programmable cluster. The motivation of using PULP is 170 twofold: the RTL and the toolchain are open-source and well 171 documented [21], and, being a MCU-based platform, PULP 172 truly represents a low-power IoT device with strict resource 173 constraints [2]. Additionally, PULP is selected to produce a 174 prototype that can be shared and adapted, without limiting 175 the compatibility of this methodology to the PULP platform, 176 as highlighted by the fact that the layer-wise approximation 177

247

strategy of Section III-B and the deployment process of 178 Section III-E do not have any hardware dependencies except 179 the multiplier. The PULP project includes libraries and tools 180 to easily export a NN model written in PyTorch to C code 181 compatible with the MCU. The MARLIN framework includes 182 modified versions of the RI5CY core and software tools, 183 adapted to include the approximation level selection through 184 control status register (CSR) instructions, whose generation 185 and compilation are added to the original PULP toolchain. 186

187 B. Approximate Neural Networks

In [19], 600 non-reconfigurable approximate multipliers are 188 tested with a multi-layer perceptron for MNIST and LeNet-5 189 for the Street View House Numbers dataset. Each approximate 190 neural network is executed using one of the 600 multipliers for 191 every convolutional layer following five retraining steps, thus 192 generating 600 different sets of weights for each model. In 193 [16], [27] is suggested that hardware-aware retraining, while 194 being a time-consuming, resource-intensive strategy, can miti-195 gate the effect of approximation. Mrazek et al. in [15] present 196 ALWANN, a framework for the approximation of NNs where 197 the assignment of each layer to an approximate multiplier, 198 among the eight-bit ones in EvoApproxLib [28], is performed 199 through the multi-objective genetic algorithm NSGA-II. The 200 parameters of each approximate NN are fine-tuned (updated 201 w.r.t. the starting model) without retraining. Therefore, for 202 each approximate configuration, a new set of weights must 203 be used for each convolutional layer. Similarly, Jain et al. in 204 [20] present a methodology to map the layers of a NN on a 205 group of systolic arrays, each composed of several instances 206 of one approximate multiplier. The arrays are part of the same 207 accelerator, with each region processing only one layer of the 208 network. Contrary to [15], the weights are not updated; there-209 fore, the original weights can be used with different approxi-210 mate configurations. However, considering several static mul-211 tiplier architectures is not a scalable approach for a general-212 purpose processor due to the area overhead; this method also 213 impacts flexibility in a custom array accelerator. Our work 214 proposes a single runtime-reconfigurable multiplier with 256 215 approximation levels, trading off the complexity of additional 216 control logic with increased flexibility. In [17], Tasoulas et al. 217 propose a methodology based on the modification of the bias 218 parameter of each layer to alleviate the approximation error. 219 Similarly to [15], this approach generates a new set of weights 220 for each approximate NN. However, their multiplier features 221 only three runtime adjustable approximation levels. Moreover, 222 the reconfiguration is handled by chaining two bits to each 223 weight stored in memory, increasing the storage requirements 224 and energy associated with data movement. According to [29] 225 and [30], the resilience of a NN model to adversarial attacks 226 depends on the approximate multiplier adopted. Consequently, 227 using different configurations, including an exact one, is sug-228 gested to achieve higher error tolerance in various scenarios. 229 MARLIN applied to [30] would remove the constraint of 230 using a single fixed approximate multiplier, enabling error 231 compensation, used in [15], [17], [20] and this work to 232 improve the accuracy. Moreover, MARLIN can eliminate the 233

limitations of [29], in which 13 different multipliers are used 234 within each PE to support 13 approximation levels, enabling 235 21.3 times more approximation levels with a thirteenth of the 236 multipliers instances. Similarly to [18], we selected a RISC-237 V-based deployment platform, thus simplifying the software 238 configuration of the approximate hardware. Nevertheless, in 239 [18], the configuration signals, handled by a control unit 240 external to the core, are generated by the user, thus relying 241 on his expertise rather than on an automated mechanism. 242 We overcome these limitations by embedding the runtime 243 approximation control inside the instructions processed by the 244 RISC-V core, leveraging the flexibility of the PULP platform. 245

III. PROPOSED METHODOLOGY





Fig. 2. MARLIN framework with hardware support specific for PULP SoC.

Our framework comprises different building blocks inter-248 acting with one another to determine a flow covering from the 249 model definition up to its hardware deployment. From now 250 on, we study the specific case of PULP platform, depicted in 251 Figure 2. The first external input required is a valid dataset, 252 such as MNIST or CIFAR-10, which defines the target task 253 of the NN model. Given a specific application, there are 254 often constraints on the minimum acceptable accuracy or the 255 maximum tolerable energy consumption. Once the training 256 dataset is provided with the specifications, a suitable NN 257 model is identified and described with PyTorch [31]. At this 258 stage, hyperparameters tuning is crucial to obtain consistent 259 results and a model that is already robust to quantization 260 errors. This phase implies choosing the number of training 26 epochs, the type of optimizer, and other learning parameters. 262 A standardized representation of the NN in the form of a 263 data flow graph is required to port the model to PULP. 264 For this reason, the trained NN is passed to NEMO [32], 265 which transforms a floating-point model to an integer one in 266 ONNX format. The precision of the model is fixed at this 267 point in the procedure, with the added constraint that the bit-268 width of weights and activations cannot be above eight bits, 269 either signed or unsigned. Up to this point of the procedure, 270 the model has no knowledge of the approximation. On the 271 hardware side, a reconfigurable inexact multiplier is designed 272

and instantiated in MARLIN with a 9x9 bits parallelism. As 273 MARLIN's optimization software only requires, for each con-274 figurable approximation level, a look-up table (LUT) storing 275 all the possible results for each couple of input operands 276 and the average power to execute a single multiplication, the 277 multiplier block is interchangeable. A single runtime recon-278 figurable approximate unit or several multipliers with fixed 279 approximation levels could be integrated into the framework 280 with little or no modification. Once this high-level description 281 of the computational unit is available, the approximate model 282 can be implemented and tested through the AdaPT library 283 [33]. Any NN topology built with PyTorch's convolutional and 284 fully connected layers can be easily included in MARLIN by 285 overloading the layer definitions with the AdaPT ones without 286 retraining or changing the model architecture. MARLIN solves 287 the complex multi-objective problem of assigning an optimal 288 approximation level to each layer through NSGA-II, described 289 in Section III-B. It requires repeated simulations of the model 290 with the selected configurations to evaluate their accuracy 291 and power consumption. The obtained Pareto front will show 292 different possible trade-offs between accuracy and power, 293 corresponding to the two fitness functions NSGA-II tries to 294 optimize. The last step that MARLIN performs on the software 295 side is the C code generation to execute the model on the target 296 hardware, presented in Section III-D. A modified version of 297 DORY [34] is used to accomplish this task. This is the first part 298 of MARLIN which requires knowledge of the actual hardware 299 architecture, including a detailed high-level description of 300 every memory level size and latency to perform memory tiling 301 effectively. For this work, PULP was selected as the target 302 platform among those supported. DORY receives as an input 303 the ONNX model generated by NEMO and an additional 304 node-by-node dictionary of the NN containing information of 305 the approximation of each layer retrieved by NSGA-II. The 306 modified DORY tool generates the C code for the provided 307 approximate architecture with the received configuration. For 308 this purpose, DORY has to be aware of the modifications the 309 PULP platform undergoes. An approximate unit is added to 310 the execution stage of the cluster cores to approximate all 311 the relevant instructions in the computation of convolutional 312 layers. It is based on the same reconfigurable multiplier, whose 313 LUTs are used by NSGA-II and AdaPT. This unit is managed 314 through a dedicated CSR in charge of activating, deactivating, 315 and configuring its approximation level. In Section III-E is 316 discussed how the modified DORY automatically inserts CSR 317 write and set instructions in the C code of the model to enable 318 the approximate unit when required. The final code runs on 319 the PULP platform through PULP-SDK. The model can be 320 321 executed by providing input data read from the external L3 memory while the weights are stored in L2 or L3 memory, 322 depending on their sizes. The support of a real hardware RTL 323 on which the model can run is crucial for the validation 324 of the proposed co-design methodology, allowing accurate 325 estimations of the metrics of interest on a complex system. 326

327 B. Genetic Search of Optimal Inter-Layer Approximation

In this work, we use NSGA-II, to solve the multi-objective problem of finding NN configurations with different trade-offs

between energy and accuracy. NSGA-II is a multi-objective 330 genetic algorithm that evolves a population of solutions using 331 non-dominated sorting and crowding distance assignment to 332 classify and rank individuals based on their dominance and 333 diversity. Crossover, i.e., recombination of different chromo-334 somes, and mutation, i.e., variation of a gene, are applied to 335 create offspring solutions, which are then integrated with the 336 parent population. The selection process favors solutions from 337 less crowded Pareto fronts and those with higher crowding 338 distances, promoting the front exploration and providing a 339 diverse set of non-dominated solutions [23]. The motivations 340 for choosing the NSGA-II algorithm are its proven effec-341 tiveness in multi-objective optimization and relative ease of 342 implementation and tuning compared to other alternatives such 343 as reinforcement learning or Bayesian optimization. NSGA-II 344 generates a set of optimized approximate NN configurations 345 and selects for each one which approximation level is more 346 suitable for each convolutional layer. Algorithm 1 details the

Algorithm 1 Approximation level selection with NSGA-II

Algorithm 1 Approximation level	selection with NSOA-II
1: $\triangleright M$ is the quantized exact model	
2: $\triangleright axx_mult$ is the approximate multiplier	
3: $\triangleright Ng$ is the number of generations	
4: $\triangleright Np$ is the population size	
5: \triangleright Pc is the crossover probability	
6: $\triangleright Pm$ is the mutation probability	
7: $\triangleright \vartheta$ is the chromosome of L elements, in range	[0, A], storing the mult. configuration
8: $\triangleright f_1(\vartheta), f_2(\vartheta)$ are the fitness functions to φ	optimize
9: $\triangleright P_n$ is the population at iteration n, with si	ze Np
10: Initialization	
11: $L \leftarrow count(M.Conv)$ \triangleright	Number of Conv. layers in the model
12: $A \leftarrow A_0$ \triangleright Number	of approximation levels for multiplier
13: $P \leftarrow P_0$ \triangleright	Initial population vector randomly set
14: $f(P_0) \leftarrow (f_1(P_0), f_2(P_0))$	Initial fitness evaluation
15: ▷ Execution	
16: for $(n = 0; n < Ng; n + +)$ do	
17: $Q_n \leftarrow \text{Tournament}(P_n, Pc, Pm)$	
18: $retrain_models(M, axx_mult, Q_n)$	
19: $f_1(Q_n) \leftarrow 1/accuracy(M, axx_max_m)$	$ult, Q_n)$
20: $f_2(Q_n) \leftarrow energy(M, axx_mult, Q_n)$	(2_n)
21: $f \leftarrow (f_1, f_2)$	
22: $R_n \leftarrow P_n + Q_n$	\triangleright Total population, size $2Np$
23: for each ϑ in R_n do	
24: $Rank(\vartheta)$	
$25: F_i \leftarrow F_i \cup \vartheta$	$\triangleright F_i$ are the fronts
26: end for	
2/: Ior each ϑ in R_n do	
28: for each ϕ_k in f do	
$\begin{array}{llllllllllllllllllllllllllllllllllll$	$istance(\vartheta, \phi_k)$
21. end for 22. Orden D based on fronts and arounding	distance
32. Order n_n based on froms and crowding 33. D_{n+1} base N_n solutions in D_n	Lindete iteration counter
34: and for	> Opdate iteration counter
35: return A	N Ontimum Pareto front is returned
55. Ituin v _{best}	p optimum r acto nont is returned

NSGA-II search flow. Each chromosome has a dimension L, 348 which is the number of layers composing the NN. The alleles 349 of each gene are encoded as an integer number between 0 350 (exact level) and A, which is the number of approximation 35 levels supported by the multiplier. Single-point crossover is 352 used to combine chromosomes while maintaining inter-layer 353 dependencies between approximate configurations, a strategy 354 used in [20] to reduce the effect of computation errors on 355 the NN accuracy without retraining. At the beginning of each 356 iteration, Np approximate NNs are retrained with 10% of 357 the training split (0.1 epoch). Then, the accuracy is evaluated 358 with the validation dataset. Contrary to previous works [15], 359 [17], the accuracy of candidate inexact NNs is not evaluated 360

immediately, but after a quick retraining with a fraction of an 361 epoch. By leveraging partial retraining and validation, each 362 NN configuration is evaluated by its resilience to computation 363 errors and the retraining effort required to recover from 364 such errors. Solutions with faster recovery will have higher 365 validation accuracy than others that are less fit, using the 366 same number of training samples, and therefore have an 367 evolutionary advantage. Therefore, retraining (or fine-tuning) 368 is used to compensate for the error and enhance the design 369 space exploration. For what concerns the fitness evaluation in 370 Algorithm 1, the inference energy is estimated as in [15], [17] 371 by multiplying the number of multiplications of each layer 372 of the model M with the average energy of the approximate 373 multiplier when set to the approximation level defined by the 374 corresponding gene of chromosome ϑ . After the evaluation 375 of the two fitness functions, the algorithm continues with the 376 mutation, crossover, tournament selection, ranking, and finally, 377 the evaluation of the crowding distance and the generation 378 of the new front. The cycle starts anew until all the Nq379 generations have been evaluated. 380

381 C. Reconfigurable Approximate Multiplier

In this work, we employ a single-cycle multiplier architec-382 ture to introduce minimum modifications in the control flow 383 of the RI5CY core. As our primary purpose is guaranteeing 384 maximum versatility and ensuring portability to different hard-385 ware with minimum effort, a parallel multiplier architecture 386 based on the Dadda reduction tree [35] is selected. The Dadda 387 algorithm is employed to compress, through half-adders and 388 full-adders, the matrix of partial products generated in the first 389 step of the multiplication, following an as late as possible 390 approach. It is preferred to the Wallace structure as it shows 391 lower delay and complexity. The modified Baugh-Wooley al-392 gorithm [36] is employed to handle signed multiplication with 393 minimum overhead in the size of the partial product matrix. 394 Both techniques are general and scalable to arbitrary operand 395 bit-width. A variation of the truncation mechanism proposed 396 in [37] is identified as the target strategy to manage the 397 dynamic setting of approximation and precision. It allows for 398 easy support of approximate and exact configurations, both for 399 full and reduced bit-width of the operands. Truncation relies 400 on a masking signal to select specific columns of the partial 401 product matrix to fix at zero to reduce the switching activity, 402 hence dynamic power, of the logic gates in that section of 403 the matrix, at the expense of an incorrect output. The number 404 of reconfiguration levels is selected considering that when the 405 approximation is extended to the most significant half of the 406 result, the error becomes unbearable, as argued in [38] and 407 [39]. An externally configurable masking signal noted as a, is 408 introduced to manage the selection of the approximation level, 409 as shown in Figure 3 for the case of 9-bit inputs and a on eight 410 bits. Contrarily to [39], each mask bit a_i of a corresponds to a 411 column of the matrix; there is no sharing of the configuration 412 signals. From bit 2 onward, the probability of the output 413 bit being one or giving a carry-out is higher than 50%, as 414 intuitively proven by the fact that the number of bits stacked 415 in the columns of the matrix is greater than two (Figure 3). 416

This high likelihood justifies the choice to gate the first row of 417 the partial product matrix, from position 2 to 7, to one, when 418 the corresponding bit of the approximation mask is active, and 419 it is implemented through the OR with a_i complemented in 420 Figure 3. All other bits are gated to zero, similarly to [37], 421 using two-inputs AND gates for the masking. An additional 422 feature of the designed architecture, which is uncommon to 423 most approximate multipliers, is the capability of runtime 424 reconfiguration of the precision of the operands. A precision 425 masking signal p is introduced to perform data-gating on 426 the partial product matrix using a mechanism similar to the 427 approximation. The p mask is externally configured according 428 to the precision of the expected result, thus allowing power 429 saving when mixed-precision multiplications are required. The 430 minimum supported bit-width for the input operands is fixed 431 to two. The precision mask has the length of the output minus 432 four. Figure 3 shows the precision signal on fourteen bits p_i 433 covering the most significant part of the partial products. When 434 a precision mask bit is set to zero the corresponding column 435 of the matrix is entirely zeroed. The complete logic inserted in 436 the generation of the partial product matrix to manage multiple 437 approximation and precision levels is depicted in Figure 3.



Fig. 3. Precision and approximation configuration management for the proposed multiplier. The approximation level is selected with the mask a, while the precision is selected with the mask p. a_j indicates the j^{th} bit of the approx_mask signal a, p_j the j^{th} bit of the precision_mask signal p, while pp_{ij} is the j^{th} bit of the i^{th} partial product evaluated according to modified Baugh-Wooley algorithm.

D. Reconfigurable Approximate RISC-V Core for PULP SoC

The RI5CY core architecture instantiated in the PULP 440 cluster must be adapted to enable reconfigurability in terms 441 of approximation and precision. The first issue to tackle to 442 introduce inexact operators in the core is how to expose them 443 in the instruction set architecture (ISA) so that a programmer 444 can effectively use them. A possible solution is the one pro-445 posed in [18]: adding custom instructions that, when decoded, 446 configure the execution stage to use approximate operators. 447 However, this approach requires an additional instruction for 448 each inexact arithmetic operation supported by the hardware, 449 defined with a new custom format capable of encoding the 450 approximation level. Another drawback is related to the fact 451 that, in this specific case, the C executable code is generated 452 automatically by the DORY tool. If custom instructions were 453 used, the user would have to replace the standard instructions 454 with the custom ones, whenever necessary, analyzing the 455

generated C code line-by-line. The same could be achieved by 456 making the compiler aware of the approximated instructions 457 and where they are needed, but that would be time-consuming 458 to implement and maintain. The methodology suggested in 459 this work addresses flexibility and simplicity by defining a 460 new custom CSR handling all the control and configuration 461 of approximate operators. This approach is scalable; a single 462 32-bit register can manage all the new operations and does 463 not occupy additional instruction encoding space. It also 464 enables reconfiguration, as part of the register bits control 465 the approximation and precision level. Finally, it is much 466 more programmer-friendly as it does not require significant 467 changes in the C code of the microcontroller, except for the 468 addition of CSR instructions. In order to achieve these results, 469 the proposed methodology to enable the online configuration 470 of the approximate multiplier relies on the following steps. 471 First, a CSR write instruction sets the precision_mask 472 and approx_mask fields according to the specification of 473 the layer. Secondly, before the computation starts, a CSR 474 set instruction enables approx_mac and approx_dot8, 475 which are disabled when the computation is over. Each CSR

Algorithm 2 PULP-NN matmul function pseudocode

1: \triangleright ch in/ch out input/output channels of the Conv layer 2: ▷ k_x/k_y filter dimension along x/y 3: \triangleright im2col is ch_in $\cdot k_x \cdot k_y$ 4: \triangleright col cnt im2col is im2col & 0x3 5: \triangleright chan_left is ch_out & 0x3 6: > Initialization > 26 instr 7: Load params. from stack and define variables 8: > Execution 9: for $(i = 0; i < ch_out >> 2; i + +)$ do $10 \cdot$ Setup ▷ if first iteration 41 instr, else 6 instr 11: for (j = 0; j < im2col >> 2; j + +) do 12: > 15 instr Setup 13: Multiply and accumulate + save > 15 instr $\cdot (im2col >> 2) + 9$ instr 14 end for 15: if (im2col >> 2 == 0) then; Setup end if ▷ 12 instr 16: while $(col_cnt_im2col ! = 0)$ do ⊳ 4 instr 17. Setup 18: Multiply and accumulate + save \triangleright 15 instr $\cdot (im2col >> 2) + 2$ instr 19. end while 20: Quantize and save results ⊳ 54 instr 21: end for 22. Setup ⊳ 10 instr 23: while (chan_left) do 24: ▷ if first iteration 23 instr, else 1 instr Setup 25: for (j = 0; j < im2col >> 2; j + +) do 26: Setup ⊳ 6 instr 27. Multiply and accumulate + save \triangleright 6 instr $\cdot (im2col >> 2)$ +4 instr 28 end for 29: if (im2col >> 2 == 0) then; Setup end if ⊳ 6 instr 30. while $(col \ cnt \ im2col ! = 0)$ do 31: Setup ⊳ 4 instr 32: Multiply and accumulate + save \triangleright 6 instr $\cdot (im2col >> 2) + 1$ instr 33: end while 34: Quantize and save results ▷ 13 instr 35: end while Save parameters and return 36: ▷ 24 instr

instruction takes one clock cycle. In the general case, the 477 last couple of CSR instructions are executed a number of 478 times which depends on the tiler split performed by DORY. 479 Their position in the code is optimized to produce minimum 480 overhead in the control flow, considering the presence of 481 MAC instructions that must produce the correct result. The 482 usage of three instructions, rather than two, is forced by 483 the specific organization of the template C files provided by 484 DORY and PULP-NN C library [22]. The CSR instruction 485 activating the approximate unit is the one executed more 486

frequently. It is located before the matmul function, whose 487 pseudocode and number of assembly instructions are depicted 488 in Algorithm 2. In a pessimistic estimation, a CSR set is 489 performed once for each matmult function call, providing a 490 quantitative measure of the reconfiguration overhead on the 491 execution time. The GCC compiler is extended to account for 492 the new approx CSR, whose fields are given in Figure 4. The



Fig. 4. The approx CSR configuration.

493 Xpulp ISA extension [22] provides additional multiply-related 494 instructions compared to the basic ones of RV32M. Some of 495 these, such as MAC and SIMD dot products, are useful for 496 NN inference. Here, the choice is to provide an approximate 497 implementation only for instructions used in convolutional 498 and linear layers. An in-depth analysis of the disassembly 499 code produced by custom convolutions and a complete NN is 500 performed to select them. The assembly instructions included 501 in these benchmarks are approximated, together with others 502 for which it is straightforward to extend support; they are all 503 listed in Table I. In this first implementation, the approximate 504 pipeline only computes 8-bit multiplications, even when the 505 issued instruction expects a 32-bit or 16-bit operation. This 506 restriction cannot cause any error assuming that NN layers 507 quantization is always on 8 bits or below, which is the 508 ordinary case. The instructions for which approximate support 509 is provided are split into three subgroups, according to Table I. 510 For each category, the approx CSR has a configuration bit; 511 when this bit is set, all the instructions belonging to that 512 group are executed in approximate mode. Besides the custom 513 CSR, a unit responsible for inexact computation is inserted 514 in the execution stage of the pipeline alongside the exact 515 multiplier unit, as shown in the high-level block diagram of 516 the approximate core in Figure 5. This design choice requires

TABLE I APPROXIMATE INSTRUCTIONS MNEMONICS

MAC	MUL		IAC MUL DO		OT8
p.mac	mul	p.mulsN	pv.dotup.b	pv.sdotup.b	
p.macsN	p.muls	p.muluN	pv.dotusp.b	pv.sdotusp.b	
p.macuN	p.mulu		pv.dotsp.b	pv.sdotsp.b	

some modifications in the decoding phase of the instruction. 518 Based on some control signals, the decoder has to activate 519 either the correct or the inexact unit. The arithmetic block 520 that is not selected for the instruction currently in the decode 521 stage does not perform any operation in the next clock cycle 522 as its inputs are not updated. Four instances of the designed 523 multiplier are allocated in the reconfigurable approximate unit, 524 as in Figure 6. They are all used in parallel to perform SIMD 525 dot products on 8 bits, while only one is activated for MUL 526 and MAC-related operations. 527

⁴⁷⁶

557

564

565

JOURNAL OF LATEX CLASS FILES, VOL. 18, NO. 9, SEPTEMBER 2020



Fig. 5. RI5CY pipeline with approximate operations support.



Fig. 6. RI5CY approximate multiplier unit. The leftmost multiplier is in charge of MUL and MAC-related operations. The shifter is for immediate instructions with the N suffix in Table I. All multipliers work in parallel to execute dot8 operations and their result is fed to a 32-bit five-operands adder. The sign extension block handles the split of the 32-bit operands into four 8-bit chunks and their sign extension according to the decoded instruction.

528 E. Approximation in PULP Toolchain

Part of the PULP toolchain must be adapted to enable run-529 time approximation. Through NEMO and DORY libraries, the 530 PULP platform offers software support to generate executable 531 C code tailored to its architecture, starting from a PyTorch 532 NN model. NEMO is a PyTorch complementary framework 533 developed as a support tool to transform an already trained, 534 full-precision NN into an integer one, performing the quanti-535 zation and calibration of the model. DORY is an open-source 536 tool for optimizing NNs mapping on PULP and other MCUs. 537 Two building blocks of DORY, the configurable templates and 538 PULP-NN back-end functions, are modified to automatically 539 add the approximate CSR instructions in proper points of the 540 C code, while the mapping optimization is unaffected. Besides 541 an ONNX graph, the modified DORY uses as input a JSON file 542 containing a layer-by-layer description of the quantization and 543 approximation of the NN. Once these parameters become part 544 of the DORY intermediate representation, they are used to fill 545 hardware-specific template files with the correct CSR setting 546 and generate the C code for the different layers. Relying on 547 a JSON dictionary guarantees flexibility, as new items can 548 be defined for each node in the network. Furthermore, it is 549

general; at this level, every type of approximate multiplier could be available, either reconfigurable or not. Moreover, as the precision information on the layer is kept separate from the approximation level, it is possible to configure the layer as exact but with reduced precision. This choice allows to save power by leveraging operations with reduced bit-width rather than with inexact computation.

IV. EXPERIMENTAL RESULTS

This section presents the computing setup used to conduct experiments to validate the proposed methodology, summarized in Figure 7. We synthesized and tested the modified RISCY with the reconfigurable multiplier and extracted relevant hardware metrics of the core and the arithmetic operator alone. Additionally, we compare our approach against stateof-the-art techniques that apply layer-wise approximation.



Fig. 7. MARLIN framework and computing setup for software simulation.

A. Multiplier Characterization

A 9-bit reconfigurable multiplier is designed according to 566 the methodology in Section III-C to compute SIMD 8-bit dot 567 products supporting all possible combinations of signed and 568 unsigned operands in RI5CY core. The multiplier is synthe-569 sized with Synopsys Design Compiler (DC) using the UMC 570 65 nm process technology library. The design is wrapped by a 571 group of input and output registers to ease the enforcement 572 of timing constraints. Similarly to [38], the target clock 573 period is set to 2 ns, making the multiplier critical path delay 574 compatible with the one it shows on the synthesized core. 575 The compile_ultra command is issued to generate the 576 gate-level netlist. Post-synthesis simulation is performed on 577 100000 random input samples to back-annotate the switching 578 activity for accurate power estimation, as in [40]. Random 579 input samples are a common approach in literature to charac-580 terize the power profile of approximate multipliers, [28], [37]. 581 Table II shows the main hardware metrics and mean-relative 582 error distance (MRED) obtained for the maximum precision 583 configuration with different approximation levels. 584

$$MRED = \frac{1}{n} \sum_{i=1}^{n} \frac{|\hat{y}_i - y_i|}{|y_i|}$$
(1)

The MRED, a commonly used metric to estimate the performance of inexact arithmetic units [38]–[40], is defined in Equation (1), where n is the number of possible combinations of input values (i.e. $n = 2^{18}$, for a 9x9 multiplier), and \hat{y}_i 588

and y_i are the i^{th} approximate and exact results, respectively. 589 Error metrics are evaluated through exhaustive simulation on 590 the entire inputs dynamic [38], [40]. Table II presents the 591 same estimations, obtained with the same constraints and 592 testing conditions, for an exact 9x9 signed multiplier described 593 behaviorally in HDL and optimized by Synopsys DC, through 594 Synopsys DesignWare (DW) library. Our multiplier has an 595 area overhead of 28%, due to the additional logic for online 596 reconfigurability. This feature does not impact on the critical 597 path, as it remains under the 2ns constraint. Finally, our 598 multiplier, in the exact configuration, saves 41.4% of power 599 compared to the one by Synopsys DW, while the highest 600 approximate level saves up to 60.1%.

TABLE II METRICS OF THE 9X9 SIGNED RECONFIGURABLE MULTIPLIER AND COMPARISON WITH AN EXACT 9X9 MULTIPLIER

Design	Area $[\mu m^2] (GE)^1$	Arrival time [ns]	Approx level	Total power $[\mu W]$	MRED
Exact	607.3 (422)	1.7	-	414.0	0
			0	241.2	0
MARLIN	822.6 (572)	1.8	127	183.0	0.07
			255	164.4	0.18
1					

¹GE is the 2-input drive-strength-one NAND gate equivalent area



Fig. 8. MRED and dynamic power variation of the proposed multiplier with the approximation level for the full-precision 9-bit inputs configuration. The approximation level on the x-axis corresponds to the 1's complement of the approx_mask value.

Although the dynamic power variation with the approximation 683 setting does not follow a monotonic trend, as can be observed 603 in Figure 8, sub-optimal configurations are retained as they 604 come at no hardware cost in area, power and latency once the 605 8-bit approx_mask signal is inserted in the multiplier logic. 606 Every level has its error distribution, which can be optimal for 607 a specific NN layer. An optimal subset of the 256 configura-608 tions might be selected offline depending on the dataset and 609 NN analysed. However, the proposed methodology aims to 610 be application-agnostic and versatile to several use-cases, thus 611 hardware and software provide support for the most generic 612 scenario. The jumps in power and MRED in Figure 8 can be 613 explained by looking at the logic in Figure 3. The MRED 614 shows higher steps when the approximation moves towards 615 columns on the left (e.g., from configuration 63 to 64, or from 616 191 to 192). On the contrary, more power is saved when more 617 columns, as far to the left as possible, are gated, reducing the 618

switching activity of the compressors in the Dadda Tree. This 619 condition in Figure 8 corresponds to approx mask values 620 $\bar{a} = 2^{j} - 1, j \in [0, 1, 2, 3, 4, 5, 6, 7, 8]$. Consequently, when 621 approx_mask is set as above, that configuration will save 622 more power than the subsequent ones. Furthermore, it is a 623 Pareto optimal point since all the successive configurations 624 imply the next columns to the left to be approximated, thus 625 increasing the MRED.

TABLE III COMPARISON BETWEEN MARLIN 8x8 SIGNED MULTIPLIER AND STATE-OF-THE-ART APPROXIMATE MULTIPLIERS.

N 1.5 15		Exact	Online	LOF	Area	Dynamic
Multiplier	Conf	conf	reconf	MRE	$[\mu m^2]$ (GE)	nower $[\mu W]$
		com	recom			power [µm]
Ha [41]	-	X	X	0.102	461.2 (321)	170
Strollo [40]	-	X	X	0.053	496.8 (345)	181.6
	1	X	X	0.031	529.2 (368)	222.3
Yang [42]	2	X	X	0.041	516.2 (359)	206.6
	3	X	X	0.069	500 (348)	190.9
	1KV6	1	X	0	541.1 (376)	167.6
Mul8a [28]	1KX5	X	X	0.089	378 (263)	116.3
Mulos [20]	1L2N	X	X	0.274	200.9 (140)	60.9
	1L12	X	X	1.347	126 (88)	36.6
	0	1	1	0	650.9 (452)	188.3
de La Guia [37]	127	1	1	0.374	650.9 (452)	130.3
	255	1	1	1.065	650.9 (452)	109.5
	0	1	1	0	652.7 (454)	188.6
MARLIN	127	1	1	0.236	652.7 (454)	130.6
	255	1	1	0.621	652.7 (454)	109.9

626 Table III compares different state-of-the-art techniques to 627 design approximate multipliers. All multipliers are 8-bit signed 628 and have been synthesized with a clock period of 2 ns and 629 characterized as for the 9x9 case. For a fair comparison, our 630 multiplier is rescaled on 8 bits. The design choice to fix some 631 bits of the partial product matrix at one rather than zero when 632 the corresponding columns are approximated improves the 633 MRED by up to 41.7% compared to [37] truncation approach, 634 with negligible area and power overhead. Compared to [40]-635 [42], all based on approximate 4-2 compressors, our multiplier 636 has lower power consumption, for corresponding MRED, 637 while still covering a much wider error dynamic. Furthermore, 638 the approximate compressors approach, even when runtime 639 error tuning is implemented [38], [39], cannot provide an exact 640 configuration. This implies the need to pair an exact multiplier 641 with the approximate one to manage operations, such as 642 control flow ones, which must produce the correct output, 643 unacceptably increasing the total area. This consideration still 644 holds for Evoapproxlib multipliers [28]. Although, according 645 to Table III, [28] show better area and power metrics com-646 pared to ours, the absence of reconfigurability precludes their 647 compatibility with MARLIN. As in [15], the implementation 648 of layer-wise approximation with [28] requires to instantiate as 649 many multipliers as the number of approximation levels, which 650 becomes inconceivable above certain values due to control, 65 area, and power overhead. In conclusion, our multiplier offers 652 several power-error trade-offs, spanning the MRED dynamic 653 of most of the 8x8 multipliers compared in Table III, standing 654 as the cheapest and most flexible solution to implement layer-655 wise approximation. Our multiplier, unlike [28], [38], [40]-656 [42], also enables runtime selection of the result bit-width to 657 add versatility and reduce power when full precision is not 658 required. For the 9-bit architecture, output precision from 18 659 bits down to 4 bits is supported through data-gating on the 660

most significant columns of the Dadda tree, with a maximum 661 power saving of around 62%. In Table IV, we evaluate the 662 trade-off of runtime precision setting with a ResNet-20 model 663 executed with and without data-gating, and with different bit-664 widths. The average full precision power is evaluated with 665 the multiplier configured to execute 9x9 bit multiplications, 666 although input operands are quantized to lower precision. The 667 average reconfigured power is evaluated when the masking 668 signal p is set to match the input operand precision.

 TABLE IV

 ResNet-20 multiplier power with and without data-gating

Activations	Weights	Avg mult pow	Avg mult pow	Absolute	Relative
precision	precision	full precision	reconfigured	accuracy	accuracy
8	8	239.15 µW	237.93 μW	91.50%	100%
6	4	227.85 μW	170.09 μW	91.43%	99.92%
4	4	225.85 µW	125.35 µW	89.87%	98.22%

669

670 B. Approximate RISC-V Core Characterization

The RI5CY core featuring the approximate extension is 671 synthesized to extract area, delay, and power estimations using 672 Synopsys DC and the UMC 65 nm library. The clock period 673 is set to 5 ns. The timing constraints were relaxed with respect 674 to the ones for the multiplier alone to also accommodate in 675 the cycle time the four-operand 32-bit adder that follows it, 676 as can be observed in Figure 6. As a matter of fact, the 677 approximate multiplier instance in the core has a delay of 678 around 1.9 ns, which makes its implementation mapping and, 679 therefore, its energy contribution consistent with the analysis 680 previously performed with a 2.0 ns clock period. The two 681 cores with exact and reconfigurable approximate operators 682 are both synthesized using the command: compile ultra 683 -no_autoungroup - no_boundary_optimization 684 -timing -gate_clock. Table V compares area and tim-685 ing values obtained. The delay constraints are satisfied by both

TABLE V Performance and area comparison for exact and approximate RI5CY and their relative multiplier units

		Exact RI5CY	Approx RI5CY
	Exact mult	14842.8 (10k)	
Area $[\mu m^2]$ (GE)	Approx mult	-	4737.2 (3k)
	Total	60621.1 (42k)	67006.8 (47k)
Timing [n e]	Exact mult	4	.45
rinnig [<i>ns</i>]	Approx mult	-	4.41

686

designs, and from the fourth column of table V it can be 687 observed that our multiplier does not interfere with the micro-688 processor critical path which remains in the exact multiplier 689 unit. The area overhead caused by the approximate unit alone 690 is 7.8%, while the total overhead, considering the additional 691 control part and the approx CSR, is 10.5%. The extra area 692 cost is mainly due to the allocation of four reconfigurable 693 multipliers in order to manage 8-bit dot products. We consider 694 this overhead acceptable as this is the first prototype of this 695 architecture; however, it could be significantly decreased if 696 our four multipliers replaced the exact ones, which is possible 697 since they also feature a non-approximate mode. 698

⁶⁹⁹ The RTL model of the RI5CY is replaced by the gate-level

netlist for all cores in the PULP cluster to enable post-synthesis simulation and power estimation on a demonstrative use-case. A simple NN model is designed with PyTorch and used as a benchmark on MNIST dataset to verify the correct behavior of the entire framework and to collect power metrics. It comprises five convolutional layers, each followed by a ReLU activation function, and a final linear layer. The entire structure is depicted in Table VI. The model is trained for 30 epochs,

TABLE VI CUSTOM NN ARCHITECTURE DESCRIPTION

Layer name	Output size	Kernel size	Output channels	# Mult
conv1	28x28	7x7	3	115224
conv2	28x28	5x5	8	470400
max pool	14x14	3x3	8	0
conv3	14x14	3x3	10	141120
conv4	14x14	3x3	16	282240
max pool	7x7	3x3	16	0
conv5	7x7	3x3	24	169344
max pool	3x3	3x3	24	0
linear	1x10	9x24	1	2160

707 with a batch size of 32, an initial learning rate of $3 \cdot 10^{-3}$, with 708 a step factor of 0.3 every 5 epochs. Stochastic gradient descent 709 with momentum 0.9 is used with a weight decay of 10^{-3} . The 710 designed NN is run on the PULP platform for a single input 711 image. For this model, the overhead of the CSR instructions 712 execution can be quantified, according to Algorithm 2, in the 713 worst case, which is the first convolutional layer, as one CSR 714 set instruction every 403 instruction (0.25%). In the best case, 715 which is the last convolution, the extra cost is 0.026%. This 716 results in a negligible performance overhead due to the CSR 717 switching, and thus of reconfiguration, on the overall pro-718 cessing. Through Siemens QuestaSim, the VCD dumps of the 719 entire core and the approximate and exact multiplying units are 720 collected. They are used as inputs for Synopsys Power Shell 721 to extract power metrics based on the actual switching activity. 722 Simulations are performed using the multipliers configurations 723 obtained from the NSGA-II run that achieve accuracy over 724 90% with no retraining. For every configuration, an example 725 for each of the ten possible categories is fed to the network, 726 meaning numbers from zero to nine for MNIST. The exact 727 post-synthesis RI5CY core is simulated with the same inputs, 728 and the power of the accurate multiplier unit is collected; all 729 results are averaged. Table VII reports, for the configurations 730 listed in the first column, the NN test accuracy and, in the 731 third column, the power of the multiplier unit (the exact one 732 for the exact configuration in the first row, the approximate 733 one for all other cases) averaged over the ten simulations. 734 The fourth column contains the relative energy saving of the 735 multiplier unit measured post-synthesis, the last column shows 736 the relative energy saving estimated at the end of the NSGA-II 737 search, as described in Section III-B. The first row of Table VII 738 contains the power estimation of the exact RI5CY multiplier 739 unit, where all operators, including the multipliers computing 740 dot8 instructions, are described behaviorally, thus leaving the 741 architecture selection to Synopsys DW. Consequently, for a 742 meaningful comparison, the reference model for the NSGA-743 II relative estimation is the average energy consumed by a 744 behavioral 9-bit multiplier synthesized by Synopsys DC, with 745

10

798

the same settings of the approximate multiplier. Although 746 the exact multiplier always executes housekeeping operations, 747 these are a negligible fraction of its overall workload when all 748 NN multiplications are mapped on it; thus, their contribution to 749 the average power is minimal. Therefore, since the number of 750 multiplications of the NN model is fixed, and so is the time 751 the multiplier stays active, the relative energy is considered 752 equivalent to the ratio between the average power of the 753 approximate unit in the modified RI5CY and that of the exact 754 one in the unmodified core. 755

TABLE VII POWER CONSUMPTION AND ENERGY SAVING OF THE MULTIPLYING UNIT WITH DIFFERENT LAYER-WISE CONFIGURATIONS FOR THE TARGET NN

	Test	Average	Relative	Relative
Configuration	Acc	RI5CY mult	RI5CY mult	NSGA-II
	[%]	power $[\mu W]$	energy saving	energy saving
Exact	98.7	10.46	0%	0%
[0, 0, 0, 0, 0]	98.7	2.606	75%	42%
[59, 31, 15, 12, 3]	98.7	2.509	76%	46%
[59, 31, 15, 31, 31]	98.5	2.474	76%	48%
[255, 63, 15, 31, 11]	97.8	2.412	77%	50%
[255, 126, 3, 63, 63]	91.3	2.403	77%	53%

The fourth column of Table VII shows that the obtained 756 average energy saving on the multiplier unit is at least 75% 757 when comparing the exact core and the approximate one with 758 all multipliers configured as accurate (first and second row 759 in Table VII). The maximum saving is 77%, which is more 760 than 20% higher than the high-level estimation performed 761 in NSGA-II. However, the advantage of adopting different 762 approximate configurations is heavily reduced compared to 763 the initial evaluation. This can be observed by rescaling the 764 energy results in the fourth and fifth columns of Table VII 765 with respect to the exact configuration of our multiplier. The 766 highest estimated energy reduction, with respect to the model 767 using the exact configuration of our multiplier for all layers, 768 is 7.7%, with an accuracy loss of 7.4%, while the predicted 769 saving was 20.1%. The cause of the gain drop, obtained by 770 configuring the multiplier with a higher approximation, has to 771 be addressed to the chosen task for the network. When the 772 average power of the multiplier is estimated, 100000 random 773 values with uniform distribution are used as inputs to the 774 multiplier. However, the MNIST dataset is composed of black 775 numbers on a vast white background, which means that the 776 network inputs and, consequently, those of the multipliers are 777 not uniformly distributed. Both input images and intermediate 778 activations show a high concentration of zeros, contrary to 779 the initial assumption. The main consequence of most values 780 being zero is that data-gating, which is the primary source 781 of power saving when approximation is applied, becomes 782 ineffective as the operands are already zeros and have a low 783 switching probability themselves. A higher sensitivity of our 784 multiplier to zero input operands, compared to the DW one, 785 could also explain the increased relative energy saving in the 786 example with respect to the estimates used during the NSGA-787 II search. Even acknowledging the difference in the statistics 788 of the operands, we run the optimization algorithm using the 789 estimated average energy computed with uniformly distributed 790 inputs, keeping it data-agnostic, as commonly done in state-of-791 the-art optimization frameworks, where performance metrics 792

are estimated independently from the statistics of the dataset 793 [3], [10]–[12]. The latter enables our search algorithm to 794 generalize to new data and thus demonstrate the efficacy of 795 our method, an approach that was also used in [15]–[17], [24] 796 to estimate the energy reduction with approximate multipliers. 797

C. Benchmark with CIFAR-10

We used 6 variations of the ResNet model architecture [43] 799 to experiment with shallow and deep NNs for image classifi-800 cation, testing the effectiveness of MARLIN with CIFAR-10 801 [44], a more challenging dataset than MNIST. We based the 802 implementation of our ResNet models on the original paper 803 and used the same model architecture and hyper-parameters, 804 with 44k iterations instead of 64k, substituting the original 805 multi-step scheduling of the learning rate with a cyclical 806 scheduler [45], ranging between 10^{-1} and 10^{-4} . We opted for 807 a cyclical learning rate instead of a stationary one to achieve 808 faster convergence with fewer training iterations during the 809 genetic search. We carried out separate quantization-aware and 810 full precision training for the INT8 and FP32 models. All the 811 experiments with approximate multipliers use the INT8 quan-812 tized models; the FP32 results are presented only to provide 813 a comparison with full precision. We used scale quantization 814 with the straight-through-estimator for the weights [13], while 815 the activations are quantized using PACT [46]. Table VIII 816 reports the NNs used in the experiments.

TABLE VIII NEURAL NETWORKS USED IN THE EXPERIMENTS

Neural	# Conv.	# Mult.	FP32	INT8	Design space
network	layers		accuracy	accuracy	SIZE
ResNet-8	7	12.2M	85.33%	85.43%	$72 \cdot 10^{15}$
ResNet-14	13	26.4M	90.17%	90.32%	20.10^{30}
ResNet-20	19	40.6M	91.77%	91.50%	57.10^{44}
ResNet-32	31	68.9M	92.65%	92.58%	$45 \cdot 10^{74}$
ResNet-50	49	111.3M	92.88%	92.60%	10.10^{117}
Resnet-56	55	125.5M	93.14%	93.11%	$28 \cdot 10^{131}$

817 The INT8 accuracy results are evaluated using the approximation level 0 of the proposed multiplier, which provides exact 819 results, for all the layers of each network. The multiplications 820 reported in Table VIII are evaluated for the inference of one 821 32x32x3 input image from the CIFAR-10 dataset. The design 822 space size reported in the rightmost column is evaluated as 823 the number of unique approximate layer-wise configurations, 824 evaluated as A_x^L , with A_x being the number of approximation 825 levels of the reconfigurable multiplier, in our case 256, and L 826 the number of layers that can be approximated. 827

For ResNet-8, ResNet-14, and ResNet-20, we ran the 828 genetic search for 80 generations with a population of 70 829 individuals, whereas for ResNet-32, ResNet-50, and ResNet-830 56, we increased the generations to 120. We set both the 831 mutation probability P_m and the crossover probability P_c 832 to 0.8. During the search phase, every approximate NN is 833 retrained with 10% of the training set and the accuracy is 834 evaluated with the validation set (5000 unseen images from the 835 training set). We did not use the test set during the search phase 836 as it would have biased the results and negatively affected the 837



Fig. 9. Top 1% accuracy and normalized energy variation with different ResNet configurations. The *Pareto valid*. blue marks represent the validation accuracy evaluated during the genetic search, whereas the *Pareto test* orange marks represent the corresponding approximate configuration tested after the final retraining.

genetic algorithm. In this phase, the NN accuracy influences 838 the evolution of each individual's configuration, i.e., the layer-839 wise approximation; therefore, to obtain NN models that can 840 generalize on new unseen data and prove the effectiveness of 841 our methodology, we removed any correlation with the final 842 test results. Finally, once the Pareto front has been computed, 843 the approximate NNs are retrained for one full epoch and 844 evaluated using the test set. 845

The accuracy and energy results of only the dominant 846 solution after the full retraining of the Pareto front are reported 847 in Figure 9, providing a visual representation of the energy-848 accuracy trade-offs our methodology offers. For every NN 849 under test, the absolute accuracy difference between Pareto 850 valid. and Pareto test marks of each configuration is usually 851 smaller than 1.5%, with even smaller values for deeper models. 852 Therefore, the validation accuracy could represent a good ap-853 proximation of the final results, justifying the choice of using 854 it in the proposed search strategy. Figure 10 depicts the utiliza-855 tion of approximation levels for each NN layer. It is possible 856 to notice the presence of peaks around levels with an index 857 equal to or smaller than $2^{j} - 1, j \in [0, 1, 2, 3, 4, 5, 6, 7, 8]$, 858 corresponding to the Pareto optimal points of Figure 8. In 859 Figure 10, it is shown how the majority of approximation 860 levels are used in the Pareto front, justifying the choice to 861 maintain power and MRED-dominated configurations as the 862 most efficient levels might not be optimal ones to achieve a 863 high task accuracy. Moreover, when concatenated appropri-864 ately, some approximation levels, even the Pareto-dominated 865 ones, might mitigate the effect of computation errors on the 866 final results, a strategy used in [20] to reduce the accuracy 867 degradation. However, Figure 10 also highlights that some 868 approximation levels are never used in this use case. Future 869 development should add the possibility of pruning the search 870 space removing unused solutions or those with the lowest uti-871 lization. In these experiments, to prove that our methodology 872 is effective with an ample search space, low- and zero-usage 873 solutions are deliberately kept to test the search algorithm in 874 a worst-case scenario with the highest complexity. 875

TABLE IX Comparison with ALWANN [15] with 0.5% and 1% relative accuracy degradation

11

	This work			ALWANN		
Neural	Absolute	Relative	Energy	Absolute	Relative	Energy
network	accuracy	accuracy	Lifergy	accuracy	accuracy	Lifergy
ResNet-8 0.5%	85.21%	99.74%	77.62%	83.16%	99.88%	84.31%
ResNet-8 1%	84.59%	99.02%	69.80%	Same solutions as ResNet-8 0.5%		
ResNet-14 0.5%	89.98%	99.63%	73.32%	85.42%	99.85%	74.34%
ResNet-14 1%	89.50%	99.09%	71.64%	84.77%	99.09%	70.85%
ResNet-50 0.5%	92.14%	99.50%	80.67%	89.08%	99.92%	78.47 %
ResNet-50 1%	91.70%	99.03%	76.67%	88.58%	99.36%	70.02 %

understand how MARLIN stands against the state-of-the-art. 877 To make a fair comparison, we compare approximate NNs 878 with weights updated after a single epoch retraining for MAR-879 LIN and weight fine-tuning for ALWANN. Our method can 880 achieve better results for shallower NNs such as ResNet-8 and 88 maintain the same relative gains for ResNet-14, whereas it was 882 not able to achieve higher energy efficiency than ALWANN for 883 ResNet-50. Comparing the absolute accuracy of the NN mod-884 els, the approximate ResNet-14 within 1% relative accuracy 885 degradation outperforms all the ResNet-50 models presented 886 by ALWANN in top-1 accuracy and, by extension, in energy 887 efficiency, as ResNet-14 has 76.3% fewer multiplications than 888 ResNet-50. Our methodology is competitive, considering that 889 the multipliers used in ALWANN have better area and power-890 MRED metrics, but are not reconfigurable [15], [28]. The 891 main advantage of a reconfigurable multiplier against several 892 arrays of fixed multipliers is that it is possible to improve 893 the energy efficiency of arithmetic operations with lower area. 894 This approach extended to a systolic array, would require a 895 single array with the same multiplier architecture, whereas 896 ALWANN requires N separate sub-arrays in order to support 897 N approximation levels. 898

Table X compares MARLIN with the results presented in
[17], without including absolute accuracy metrics, as they
are not reported. We consider approximate NNs with one-
epoch retraining for MARLIN, and approximate NNs with
weight fine-tuning with and without additional bias for [17].
Compared to the NNs with no additional bias, the approximate
NNs configurations found with MARLIN require up to 13.1%990
900

Table IX compares our approach with ALWANN [15] to



Fig. 10. Approximation levels utilization for all the configurations found for each ResNet model.

less energy for ResNet-20, up to 13% less energy for ResNet-906 32, and up to 15.1% for ResNet-56. When an additional error 907 correction bias is added to the convolutional layer in [17], 908 MARLIN can still achieve up to 9,8% less energy for ResNet-909 20, 8.6% for ResNet-32, and 7,3% for ResNet-56, without 910 increasing the number of parameters and operations. Using 911 more approximation levels proved to be an effective way to 912 further reduce the inference energy, as we had 256 configura-913 tions against the 3 used in [17]. We justify these results with 914 the traditional weight-update strategy used in this work and the 915 presence of more approximate configurations. Retraining each 916 configuration is slower and more computationally expensive 917 than the multiplier-specific fine-tuning of [17], but allows a 918 better adjustment of the NN parameters to compensate the 919 computation errors, resulting in higher accuracy.

 TABLE X

 Normalized energy comparison with [17] with 0.5%, 1%, and 2%

 Accuracy degradation

	Normalized energy				
Neural	Ours	[17]	[17]		
network	Ours	w/o bias	with bias		
ResNet-20 0.5%	75.91%	86.1%	83.1%		
ResNet-20 1%	74.46%	85.6%	83.0%		
ResNet-20 2%	74.46%	85.1%	82.5%		
ResNet-32 0.5%	77.21%	85.7%	81.7%		
ResNet-32 1%	74.50%	85.7%	81.7%		
ResNet-32 2%	74.39%	85.5%	81.4%		
ResNet-56 0.5%	79.87%	94.0%	83.0%		
ResNet-56 1%	77.12%	86.1%	83.0%		
ResNet-56 2%	77.04%	86.1%	83.0%		

920

921 D. Compatibility with Other Accelerator Architectures

Two important modifications are necessary to port MARLIN 922 to other platforms: adapting the hardware architecture and the 923 mapper to include the multiplier and the configuration instruc-924 tions. The former would require an additional 8-bit control 925 signal from the PE control unit to set the approximation level. 926 The elongated critical path delay due to the approximation 927 logic could be a problem for some accelerators, but it is not 928 for [8], [9], which have a critical path compatible with the 929 proposed multiplier. For what concerns the mapper, recalling 930 the discussion of Section II-A, a modification similar to what 931 has been done in Section III-E can be implemented, inserting 932 custom instructions to configure the multiplier, with negligible 933 impact on the execution time. Since the scheduling does not 934

change, the number of computation cycles would also be 935 unaffected. Table XI reports the area and the energy overheads 936 of including and controlling the approximate multipliers in 937 three accelerators, mapping the 19 convolutional layers of 938 the ResNet-20 with 1% accuracy degradation and 74.46% 939 energy of Table X. We used the power model in Timeloop 940 [10] to evaluate the energy used to communicate to the PEs 941 the approximation level of each layer, assuming one off-chip 942 to on-chip memory transfer, and then #PEs transfers from 943 the on-chip memory to the PEs' registers. The configuration 944 energy of the entire NN is always below 0.002% of the total 945 energy evaluated with Timeloop. The area overhead of 35% 946 against exact multipliers can be negligible, considering that 947 they account for less than 10% of the PE area in [8], [9]. 948

12

949

TABLE XI MARLIN'S AREA AND COMMUNICATION OVERHEAD APPLIED TO OTHER HW ACCELERATORS

	Eyeriss [9]	DianNao [7]	Simba [8]
# PEs	256	256	1024
PE conf. comm. energy [pJ]	2138	1997	2369
(relative)	0.001%	0.001%	0.002%
Mult. area exact $[\mu m^2]$ (GE)	155468 (108k)	155468 (108k)	621875 (432k)
Mult. area approx. $[\mu m^2]$ (GE)	210585 (146k)	210585 (146k)	842342 (586k)
(relative)	(+35%)	(+35%)	(+35%)

E. Discussion

The optimization approach of Section III-B allowed MAR-950 LIN to outperform previous works that relied on parameter 951 fine-tuning [15], [17], leveraging partial retraining. MARLIN 952 was run on a 32-thread Ryzen 5950X CPU with 64GB DDR4 953 RAM and an Nvidia Quadro RTX A5000. The GPU was 954 used only during the initial training of the FP32 and exact 955 INT8 NNs presented in Table VIII, while the CPU was 956 used to simulate the approximate convolutional layers during 957 the training, validation, and test done during the search, as 958 AdaPT only supports CPU computation [33]. The number of 959 threads used during the computation was set to 16 for every 960 experiment to compare how MARLIN execution time scales 961 with different NNs depths. Table XII reports the execution time 962 for the search phase and the training of the last Pareto front of 963 Figure 9. On average, partial retraining is $\approx 6x$ faster than full 964 retraining. An alternative implementation that leverages the 965 GPU processing power, based on [47], is in development, with 966 This article has been accepted for publication in IEEE Transactions on Circuits and Systems--I: Regular Papers. This is the author's version which has not been fully edited and content may change prior to final publication. Citation information: DOI 10.1109/TCSI.2024.3365952

JOURNAL OF LATEX CLASS FILES, VOL. 18, NO. 9, SEPTEMBER 2020

the objective of reducing the search time with more complex 967

NN models, allowing for a broader search space analysis. 968

Compared to [15], the iteration time during the search phase 969

is reduced by 72.8% for ResNet-8, 92.3% for ResNet-14, and 970

85.4% for ResNet-50. This speed-up is due to the increased 971

utilization of CPU threads, as our training loop processes more 972 images during each iteration compared to [15].

TABLE XII MARLIN'S AVERAGE EXECUTION TIME WITH 16 THREADS

Neural network	Search phase		Final training	
	One iter.	Total	One iter.	Total
ResNet-8	6.8 sec.	10.6 hours	40.9 sec.	24.6 min.
ResNet-14	7.7 sec.	12 hours	79.4 sec.	35.7 min.
ResNet-20	19.6 sec.	30.5 hours	115.7 sec.	90.6 min.
ResNet-32	30.3 sec.	70.7 hours	189.0 sec.	81.9 min.
ResNet-50	47.1 sec.	109.9 hours	296.2 sec.	69.1 min.
ResNet-56	56.9 sec.	132.8 hours	329.2 sec.	93.3 min.

973

A limitation in finding the optimal trade-off between energy 974 and accuracy is the dimension of the search space, which 975 determines the search time and requires a carefully tuned 976 search strategy. This problem is also found in mixed precision 977 layer-wise quantization, in which the search space is q^{2L} , 978 with q quantization levels for weights and activations, for L979 layers [11], [12]. Future work should focus on pruning the 980 search space after a number of experiments (i.e., NSGA-II 981 generations) to reduce its size, possibly reducing the time to 982 converge. A further improvement over layer-wise approxima-983 tion can be proposed by looking at past works on quantization. 984 In AutoQ [3] channel-wise quantization is used to reduce the 985 inference energy with less accuracy degradation than [11], 986 [12], proving that NNs resilience to quantization errors has an 987 intra-layer dependency besides the inter-layer one. Therefore, 988 achieving an optimal energy-accuracy trade-off is possible by 989 extending approximate computing in the channel dimension. 990 In AutoQ [3], a bit-serial accelerator is required to support 991 channel-wise quantization, whereas with MARLIN the only 992 necessary modification, to enable it with the proposed RISC-993 V core, would be to adapt the CSR instructions inserted by our 994 modified version of DORY. Nonetheless, the main challenge 995 would be the efficient exploration of a wider search space. 996

997

V. CONCLUSION

In this paper, we presented MARLIN, a layer-wise approxi-998 mation methodology leveraging a single multiplier architecture 999 that can be configured runtime with 256 approximation levels 1000 to achieve an optimal trade-off between the inference energy 1001 and the task accuracy. In this work, a prototype based on a 1002 modified RI5CY core is proposed to test our methodology 1003 on a low-power IoT platform. The PULP toolchain has been 1004 adapted to automatically include the runtime approximation 1005 level selection alongside the instructions executed while pro-1006 cessing convolutional layers. MARLIN can evaluate thousands 1007 of different NNs, leveraging NSGA-II to find the optimal 1008 configuration by generating a Pareto front that contains a set 1009 of layer-wise approximate NNs with reduced inference energy, 1010 without a significant accuracy loss. 1011

REFERENCES

[1] M. Shafique, T. Theocharides, V. J. Reddy, and B. Murmann, "TinyML: 1013 Current progress, research challenges, and future roadmap," in 58th 1014 ACM/IEEE Design Automation Conf. (DAC), pp. 1303-1306, 2021. 1015

13

1012

1051

1052

1053

- [2] M. Merenda, C. Porcaro, and D. Iero, "Edge machine learning for AI-1016 enabled IoT devices: A review," Sensors, vol. 20, no. 9, p. 2533, 2020. 1017
- [3] Q. Lou, F. Guo, M. Kim, L. Liu, and L. Jiang, "AutoQ: Automated 1018 kernel-wise neural network quantization," in Intl. Conf. on Learning 1019 Representations, 2019. 1020
- N. Fasfous et al., "AnaCoNGA: Analytical HW-CNN co-design using [4] 1021 nested genetic algorithms," in Design, Automation & Test in Europe 1022 Conference & Exhibition (DATE), pp. 238-243, 2022. 1023
- [5] A. Garofalo, M. Rusci, F. Conti, D. Rossi, and L. Benini, "Pulp-nn: 1024 accelerating quantized neural networks on parallel ultra-low-power risc-1025 v processors," Philosophical Trans. of the Royal Society A, vol. 378, 1026 no. 2164, p. 20190155, 2020. 1027
- [6] A. Garofalo, G. Tagliavini, F. Conti, D. Rossi, and L. Benini, "XpulpNN: 1028 Accelerating quantized neural networks on RISC-V processors through 1029 ISA extensions," in Design, Automation & Test in Europe Conference 1030 & Exhibition (DATE), pp. 186-191, 2020. 1031
- T. Chen et al., "Diannao: A small-footprint high-throughput accelerator [7] 1032 for ubiquitous machine-learning," in *Proceedings of the 19th Intl. Conf.* 1033 on Architectural Support for Programming Languages and Operating 1034 Systems, p. 269-284, 2014. 1035
- [8] Y. S. Shao et al., "Simba: Scaling deep-learning inference with multi-1036 chip-module-based architecture," in Proceedings of the 52nd Annual 1037 IEEE/ACM Intl. Symposium on Microarchitecture, p. 14-27, 2019. 1038
- [9] Y.-H. Chen, T. Krishna, J. S. Emer, and V. Sze, "Eyeriss: An energy-1039 efficient reconfigurable accelerator for deep convolutional neural net-1040 works," IEEE Journal of Solid-State Circuits, vol. 52, no. 1, pp. 127-1041 138, 2017. 1042
- [10] A. Parashar et al., "Timeloop: A systematic approach to dnn accelerator 1043 evaluation," in IEEE Intl. Symposium on Performance Analysis of 1044 Systems and Software (ISPASS), pp. 304-315, 2019. 1045
- [11] K. Wang, Z. Liu, Y. Lin, J. Lin, and S. Han, "HAQ: Hardware-aware 1046 automated quantization with mixed precision," in IEEE/CVF Conf. on 1047 Computer Vision and Pattern Recognition, pp. 8604-8612, 2019. 1048
- [12] N. Fasfous et al., "HW-FlowQ: A multi-abstraction level HW-CNN co-1049 design quantization methodology," ACM Trans. Embed. Comput. Syst., 1050 vol. 20, sep 2021.
- [13] H. Wu, P. Judd, X. Zhang, M. Isaev, and P. Micikevicius, "Integer quantization for deep learning inference: Principles and empirical evaluation," CoRR, vol. abs/2004.09602, 2020.
- [14] H. Jiang, F. J. H. Santiago, H. Mo, L. Liu, and J. Han, "Approximate 1055 arithmetic circuits: A survey, characterization, and recent applications," 1056 Proceedings of the IEEE, vol. 108, no. 12, pp. 2108-2135, 2020. 1057
- [15] V. Mrazek, Z. Vasicek, L. Sekanina, M. A. Hanif, and M. Shafique, 1058 "ALWANN: Automatic layer-wise approximation of deep neural network 1059 accelerators without retraining," in IEEE/ACM Intl. Conf. on Computer-1060 Aided Design (ICCAD), 2019. 1061
- [16] X. He, W. Lu, G. Yan, and X. Zhang, "Joint design of training 1062 and hardware towards efficient and accuracy-scalable neural network 1063 inference," IEEE Journal on Emerging and Selected Topics in Circuits 1064 and Systems, vol. 8, no. 4, pp. 810-821, 2018. 1065
- [17] Z.-G. Tasoulas, G. Zervakis, I. Anagnostopoulos, H. Amrouch, and 1066 J. Henkel, "Weight-oriented approximation for energy-efficient neural 1067 network inference accelerators," IEEE Trans. on Circuits and Systems 1068 I: Regular Papers, vol. 67, no. 12, pp. 4670-4683, 2020 1069
- [18] I. Taştan, M. Karaca, and A. Yurdakul, "Approximate CPU design for 1070 iot end-devices with learning capabilities," Electronics, vol. 9, no. 1, 1071 2020.1072
- [19] M. S. Ansari, V. Mrazek, B. F. Cockburn, L. Sekanina, Z. Vasicek, 1073 and J. Han, "Improving the accuracy and hardware efficiency of neural 1074 networks using approximate multipliers," IEEE Trans. on Very Large 1075 Scale Integration (VLSI) Systems, vol. 28, no. 2, pp. 317-328, 2020. 1076
- [20] P. Jain, S. Huda, M. Maas, J. E. Gonzalez, I. Stoical, and A. Mirhoseini, 1077 "Learning to design accurate deep learning accelerators with inaccurate 1078 multipliers," in Design, Automation & Test in Europe Conference & 1079 Exhibition (DATE), pp. 184-189, 2022. 1080
- [21] PULP platform. Accessed: Jan. 23, 2023 [Online]. Available: 1081 https://pulp-platform.org. 1082
- M. Gautschi et al., "Near-threshold RISC-V core with DSP extensions [22] 1083 for scalable IoT endpoint devices," IEEE Trans. on Very Large Scale 1084 Integration (VLSI) Systems, vol. 25, no. 10, pp. 2700-2713, 2017. 1085

- [23] K. Deb, A. Pratap, S. Agarwal, and T. Meyarivan, "A fast and elitist mul-1086 tiobjective genetic algorithm: NSGA-II," IEEE Trans. on Evolutionary 1087 Computation, vol. 6, no. 2, pp. 182-197, 2002. 1088
- [24] M. Pinos, V. Mrazek, and L. Sekanina, "Prediction of inference energy 1089 on cnn accelerators supporting approximate circuits," in 26th Intl. Svm-1090 1091 posium on Design and Diagnostics of Electronic Circuits and Systems (DDECS), pp. 45-50, 2023. 1092
- E. Manor and S. Greenberg, "Using hw/sw codesign for deep neural 1093 [25] network hardware accelerator targeting low-resources embedded pro-1094 1095 cessors," IEEE Access, vol. 10, pp. 22274-22287, 2022.
- [26] E. Manor and S. Greenberg, "Custom hardware inference accelerator for 1096 tensorflow lite for microcontrollers," IEEE Access, vol. 10, pp. 73484-1097 73493, 2022. 1098
- [27] S. Venkataramani, A. Ranjan, K. Roy, and A. Raghunathan, "AxNN: 1099 Energy-efficient neuromorphic systems using approximate computing," 1100 in IEEE/ACM Intl. Symposium on Low Power Electronics and Design 1101 (ISLPED), pp. 27-32, 2014. 1102
- V. Mrazek, L. Sekanina, and Z. Vasicek, "Libraries of approximate 1103 [28] circuits: Automated design and application in CNN accelerators," IEEE 1104 Journal on Emerging and Selected Topics in Circuits and Systems, 1105 vol. 10, no. 4, pp. 406-418, 2020. 1106
- 1107 [29] E. Atoofian, "Increasing robustness against adversarial attacks through ensemble of approximate multipliers," in IEEE Intl. Conf. on Network-1108 ing, Architecture and Storage (NAS), pp. 1-8, 2022. 1109
- [30] A. Siddique and K. A. Hoque, "Is approximation universally defen-1110 sive against adversarial attacks in deep neural networks?," in De-1111 1112 sign, Automation & Test in Europe Conference & Exhibition (DATE), p. 364-369, 2022. 1113
- [31] PyTorch. Accessed: Mar. 11. 2023 [Online]. Available: 1114 1115 https://pytorch.org/.
- F. Conti, "Technical report: NEMO DNN quantization for deployment [32] 1116 model," 2020. 1117
- D. Danopoulos, G. Zervakis, K. Siozios, D. Soudris, and J. Henkel, 1118 [33] 1119 "AdaPT: Fast emulation of approximate DNN accelerators in PyTorch," IEEE Trans. on Computer-Aided Design of Integrated Circuits and 1120 Systems, pp. 1-1, 2022. 1121
- A. Burrello, A. Garofalo, N. Bruschi, G. Tagliavini, D. Rossi, and [34] 1122 F. Conti, "DORY: Automatic end-to-end deployment of real-world 1123 1124 DNNs on low-cost IoT MCUs," IEEE Trans. on Computers, vol. 70, no. 8, pp. 1253-1268, 2021. 1125
- [35] L. Dadda, "Some schemes for parallel multipliers," Alta Frequenza, 1126 vol. 34, pp. 349-356, 1965. 1127
- 1128 [36] C. Baugh and B. Wooley, "A two's complement parallel array multiplication algorithm," IEEE Trans. on Computers, vol. C-22, no. 12, 1129 pp. 1045-1047, 1973. 1130
- M. de la Guia Solaz, W. Han, and R. Conway, "A flexible low power 1131 [37] DSP with a programmable truncated multiplier," IEEE Trans. on Circuits 1132 and Systems I: Regular Papers, vol. 59, no. 11, pp. 2555-2568, 2012. 1133
- [38] T. Yang, T. Ukezono, and T. Sato, "A low-power high-speed accuracy-1134 controllable approximate multiplier design," in 23rd Asia and South 1135 Pacific Design Automation Conf. (ASP-DAC), pp. 605-610, 2018. 1136
- F.-Y. Gu, I.-C. Lin, and J.-W. Lin, "A low-power and high-accuracy 1137 [39] approximate multiplier with reconfigurable truncation," IEEE Access, 1138 vol. 10, pp. 60447-60458, 2022. 1139
- A. G. M. Strollo, E. Napoli, D. De Caro, N. Petra, and G. D. Meo, 1140 [40] "Comparison and extension of approximate 4-2 compressors for low-1141 power approximate multipliers," IEEE Trans. on Circuits and Systems 1142 I: Regular Papers, vol. 67, no. 9, pp. 3021-3034, 2020. 1143
- [41] M. Ha and S. Lee, "Multipliers with approximate 4-2 compressors and 1144 error recovery modules," IEEE Embedded Systems Letters, vol. 10, no. 1, 1145 pp. 6-9, 2018. 1146
- Z. Yang, J. Han, and F. Lombardi, "Approximate compressors for error-[42] 1147 resilient multiplier design," in IEEE Intl. Symposium on Defect and Fault 1148 Tolerance in VLSI and Nanotechnology Systems, pp. 183-186, 2015. 1149
- [43] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image 1150 recognition," in Proceedings of the IEEE Conf. on computer vision and 1151 1152 pattern recognition, pp. 770-778, 2016.
- A. Krizhevsky, G. Hinton, et al., "Learning multiple layers of features [44] 1153 1154 from tiny images," University of Toronto, 2009.
- L. N. Smith, "Cyclical learning rates for training neural networks," in 1155 [45] 1156 IEEE winter Conf. on applications of computer vision, pp. 464-472, 2017.1157
- J. Choi, Z. Wang, S. Venkataramani, P. I.-J. Chuang, V. Srinivasan, 1158 [46] and K. Gopalakrishnan, "PACT: Parameterized clipping activation for 1159 quantized neural networks," 2018. 1160

[47] E. Trommer, B. Waschneck, and A. Kumar, "High-throughput approx-1161 imate multiplication models in PyTorch," in 26th Intl. Symposium on 1162 Design and Diagnostics of Electronic Circuits and Systems (DDECS), 1163 pp. 79-82, 2023. 1164





Flavia Guella received the B.S. and M.S. (both 1165 with summa cum laude) in electronics engineering 1166 from Università degli Studi di Palermo in 2020, and 1167 Politecnico di Torino in 2023, respectively. She is 1168 currently pursuing the Ph.D. program in Electronics 1169 and Communications Engineering at Politecnico di 1170 Torino, under the supervision of Prof. Maurizio 1171 Martina and Guido Masera. Her research interests 1172 include co-design methodologies for the efficient de-1173 ployment of neural networks on low-power systems 1174 and approximate computing. 1175

Emanuele Valpreda received the B.S. and M.S. 1176 (summa cum laude) degrees in electronic and com-1177 munications engineering from Politecnico di Torino, 1178 in 2017 and 2019 respectively. Now he is a Ph.D. 1179 student with the Electronics and Telecommunica-1180 tions Department on energy efficient neural network 1181 acceleration for edge computing, under the supervi-1182 sion of Prof. Maurizio Martina and Guido Masera. 1183 His current research interests include neural network 1184 compression, reliability and approximate computing. 1185 1186



Michele Caon is a Ph.D. student with the Electron-1187 ics and Telecommunications Department, Politec-1188 nico di Torino, under the supervision of Prof. Mau-1189 rizio Martina and Guido Masera. He received his 1190 B.S. and M.S. (summa cum laude) at Politecnico di 1191 Torino in 2017 and 2019. His research interests are 1192 innovative digital, integrated, programmable com-1193 puting circuits and systems. He is currently working 1194 on near-memory computing circuits embedded in 1195 heterogeneous systems-on-chip. 1196 1197



Guido Masera (SM'07) received the Dr.-Ing. 1198 (summa cum laude) and Ph.D. degrees in electronic 1199 engineering from Politecnico di Torino, Italy. He is a 1200 Professor with the Electronics and Telecommunica-1201 tions Department, Politecnico di Torino, since 1992. 1202 His research interests include several aspects in the 1203 design of digital integrated circuits and systems, with 1204 a special emphasis on high-performance architec-1205 tures for communications, forward error correction, 1206 image and video coding, cryptography and hardware 1207 accelerators for machine learning. He has more than 1208

200 publications, two patents and was a designer of several ASIC components. 1209 Dr. Masera is an Associate Editor of MDPI Electronics and a former Associate 1210 Editor of the IEEE Transactions on Circuits and Systems I, IEEE Transactions 1211 on Circuits and Systems II and the IET Circuits, Devices & Systems. 1212



Maurizio Martina received the Dr.-Ing. and Ph.D. 1213 degrees in electronic engineering and electronic and 1214 communications engineering from Politecnico di 1215 Torino, Italy, in 2000 and 2004, respectively. He is a Professor with the Electronics and Telecommunications Department, Politecnico di Torino, since 2014. His research interests include computer architecture and VLSI design of digital integrated circuits for image and video coding, forward error correction, cryptography and artificial intelligence. He has more than 100 publications and holds two patents. He 1223

served as an Associate Editor of the IEEE Transactions on Circuits and 1224 Systems-I and as a Guest Editor of several special issues, including BioCAS 1225 2017 special issue in IEEE Transactions of Biomedical Circuits and Systems 1226 and ISCAS 2023 special issue in IEEE Transactions on Circuits and Systems-1227 II. He has been part of the organizing and technical committee of several IEEE 1228 conferences, including BioCAS 2017, AICAS 2020 and PRIME 2023. 1229