

# Deploy: Enabling Energy-Efficient Deployment of Small Language Models on Heterogeneous Microcontrollers

Moritz Scherer<sup>1</sup>, Graduate Student Member, IEEE, Luka Macan<sup>2</sup>, Graduate Student Member, IEEE, Victor J. B. Jung<sup>1</sup>, Graduate Student Member, IEEE, Philip Wiese<sup>1</sup>, Graduate Student Member, IEEE, Luca Bompani<sup>1</sup>, Graduate Student Member, IEEE, Alessio Burrello<sup>1</sup>, Member, IEEE, Francesco Conti<sup>1</sup>, Member, IEEE, and Luca Benini<sup>1</sup>, Fellow, IEEE

## I. INTRODUCTION

24

**Abstract**—With the rise of embodied foundation models (EFMs), most notably small language models (SLMs), adapting Transformers for the edge applications has become a very active field of research. However, achieving the end-to-end deployment of SLMs on the microcontroller (MCU)-class chips without high-bandwidth off-chip main memory access is still an open challenge. In this article, we demonstrate high efficiency end-to-end SLM deployment on a multicore RISC-V (RV32) MCU augmented with ML instruction extensions and a hardware neural processing unit (NPU). To automate the exploration of the constrained, multidimensional memory versus computation tradeoffs involved in the aggressive SLM deployment on the heterogeneous (multicore+NPU) resources, we introduce Deploy, a novel deep neural network (DNN) compiler, which generates highly optimized C code requiring minimal runtime support. We demonstrate that Deploy generates the end-to-end code for executing SLMs, fully exploiting the RV32 cores’ instruction extensions and the NPU. We achieve leading-edge energy and throughput of 490  $\mu$ J per token, at 340 token per second for an SLM trained on the TinyStories dataset, running for the first time on an MCU-class device without the external memory.

**Index Terms**—Accelerators, compilers, embodied AI, foundation models (FMs), neural networks, TinyML.

Manuscript received 5 August 2024; accepted 10 August 2024. This work was supported in part by the Spoke 1 on Future HPC of the Italian Research Center on High-Performance Computing, Big Data and Quantum Computing (ICSC) funded by MUR Mission 4-Next Generation EU; in part by the Chips Joint Undertaking (Chips-JU) TRISTAN Project under Agreement 101095947; in part by the CONVOLVE Project under Agreement 101070374; and in part by the EU Horizon Europe project NeuroSoC under Grant 101070634. Chips-JU, CONVOLVE, and NeuroSoC receive support from the European Union’s Horizon Europe Research and Innovation Program. This article was presented at the International Conference on Compilers, Architectures, and Synthesis for Embedded Systems (CASES) 2024 and appeared as part of the ESWEK-TCAD Special Issue. This article was recommended by Associate Editor S. Dailey. (Corresponding author: Moritz Scherer.)

Moritz Scherer, Victor J. B. Jung, and Philip Wiese are with the Integrated Systems Laboratory, ETH Zürich, 8092 Zürich, Switzerland (e-mail: schermo@iis.ee.ethz.ch; jungvi@iis.ee.ethz.ch; wiesep@iis.ee.ethz.ch).

Luka Macan, Luca Bompani, and Francesco Conti are with the Department of Electrical, Electronic, and Information Engineering, University of Bologna, 40126 Bologna, Italy (e-mail: luka.macan@unibo.it; luca.bompani5@unibo.it; f.conti@unibo.it).

Alessio Burrello is with the Department of Control and Computer Engineering, Politecnico di Torino, 10129 Turin, Italy (e-mail: alessio.burrello@unibo.it).

Luca Benini is with the Integrated Systems Laboratory, ETH Zürich, 8092 Zürich, Switzerland, and also with the Department of Electrical, Electronic, and Information Engineering, University of Bologna, 40126 Bologna, Italy (e-mail: lbenini@iis.ee.ethz.ch).

Digital Object Identifier 10.1109/TCAD.2024.3443718

THE LATEST evolutions in mainstream artificial intelligence (AI) have been driven by Transformers, which have taken over from recurrent neural networks (RNNs) and convolutional neural networks (CNNs) as the leading edge models for language processing and multimodal applications [1], [2]. The success of Transformers can be primarily attributed to the emergence of the foundation model (FM) paradigm: large Transformer models extensively pretrained on the datasets spanning trillions of tokens and then fine tuned with a much lower volume of labeled data to solve the domain-specific problems. Following the success of FMs in natural language processing (NLP) [1], [3], an increasing number of fields are starting to formulate and adapt FMs for high dimensional sensor data that has traditionally been challenging to process, like decoding the neural data [4], [5], or training embodied AI agents [6], [7], which may incorporate the multimodal sensor inputs.

Operating directly on the sensory data and in a cyber-physical loop may lead to solving many outstanding challenges in fields, such as brain-machine interfaces [5] and miniaturized robotics [7]. However, to materialize this promise, models of this class need to be *embodied* in physical devices as embodied foundation models (EFMs), and they must cope with the strict constraints in terms of compute throughput, power consumption, and footprint typical of the edge devices. Unlike the data center-scale systems, which collect and aggregate sensor data over the shared resources for high-throughput processing, embodied AI systems must process the sensor data with extremely low latency and memory capacity under tight power constraints. This is particularly challenging for the smallest class of AI-oriented computers: so-called “*TinyML*” devices operating at the extreme edge, based on microcontroller (MCU)-class devices without complex operating systems or memory-management units (MMUs), relying on the user level software to implement low-level hardware management functionalities. Despite many recent successes with the previous generation deep neural networks (DNNs), the emergence of the TinyML paradigm for the EFMs faces the dual challenge of reducing FMs to a manageable size and enabling their deployment on the tiny devices.

A first concrete step in this direction is the recent introduction of small language models (SLMs): FMs with tens to a few hundred million, rather than several billion parameters [8], [9]. While most currently available FMs are focused on processing natural language at a proof-of-concept scale, the effort toward the embedded multimodal sensor inputs with small-scale, application-specific FMs offers a highly promising path for the development of this novel class of models. Much like what happened with the initial emergence of deep learning [10], the evolution of advanced TinyML applications based on the EFMs is currently prevented by the lack of suitable targets for the deployment of these models and, even more, of the deployment frameworks that enable utilizing the existing specialized hardware to its full capabilities.

Deploying tiny EFMs requires overcoming several challenges specific to the TinyML domain. Large-scale AI inference systems typically employ the heterogeneous computer architectures composed by a conventional host (e.g., an x86 processor) and a very large throughput-oriented accelerator (e.g., H100 [11], TPU [12]), which is fully exploited only at large batch sizes. Conversely, TinyML is used for latency-sensitive applications focusing on the real-time inference without batching. As a consequence, TinyML AI inference typically employs much more specialized accelerator architectures [13], [14], leading to more complex mapping and optimization challenges for the DNN deployment. Furthermore, TinyML's strict constraints on energy efficiency and MCU-class computer architecture typically require platform-specific optimization, including memory-aware tiling, static memory allocation, and latency-hiding direct memory access (DMA) scheduling, which require advanced compiler support to scale to complex DNNs like FMs. While several compilers have limited support for user-defined kernels [15], [16], configuring and extending them requires expert knowledge, and their top-down compilation approach often clashes with loosely coupled accelerators. Moreover, mainstream compilers do not address the strict memory constraints in extreme-edge devices.

In this article, we aim to remove the first barrier toward developing EFM suited for deployment on TinyML platforms: the lack of deployment frameworks that enable their efficient execution. We demonstrate, to the best of our knowledge, the first end-to-end tool flow to deploy EFMs on the heterogeneous MCU-class systems. Specifically, we demonstrate the end-to-end deployment of a TinyStories-class [8] network on *Siracusa*, an advanced MCU in TSMC 16-nm technology featuring embedded nonvolatile memory [magnetoresistive random access memory (MRAM)] and two heterogeneous compute engines, namely an octa-core RV32 compute cluster with instruction extensions for machine learning (ML) and a multimode CNN neural processing unit (NPU), N-Eureka [13]. We present the tooling and algorithms integrated within our deployment framework, Deeploy.

The contributions of this article are as follows.

- 1) We describe *Deeploy*, a customizable, domain-specific compiler designed for generating bare metal code fitting the memory constraints of extreme edge devices. Deeploy supports all the key computational primitives

needed for the execution of Transformer-based EFMs on heterogeneous extreme edge system-on-chip (SoCs) through its bottom-up compilation approach, which allows applying advanced code optimization on expert-optimized kernel templates. We further introduce a novel algorithm for solving the tiling and static memory allocation problems for multilevel software-managed caches and its integration into Deeploy.<sup>1</sup>

- 2) We benchmark common Transformer configurations, demonstrating that code generated by Deeploy maximizes engine utilization in heterogeneous, multiaccelerator SoCs. We achieve data marshaling overheads of just 9% for large workloads with high arithmetic intensity executing on the cluster cores and NPU collaboratively thanks to efficient data movement acceleration and low-overhead offloading mechanisms.
- 3) As a concrete large-scale end-to-end use-case of Deeploy and its adaptability to heterogeneous hardware platforms, we demonstrate for the first time the deployment of a TinyStories-class SLM on *Siracusa*, a state-of-the-art heterogeneous MCU. While using on-chip memory only, we achieve a throughput of 340 token per second at an energy cost of 490  $\mu\text{J}$  for autoregressive inference. We show that using the flexible deployment flow enabled by Deeploy for the same SLM allows us to implement multilayer KV caching using on-chip memory only, improving token throughput by 26 $\times$  compared to inference without caches.

The remainder of this article is organized as follows. In Section II, the previous work in quantized neural networks, SLMs, and neural network deployment for the extreme edge devices is introduced and discussed. Section III introduces the Deeploy and discusses its deployment flow for the Transformers. Section IV discusses the SLM architecture used in this work and the approach to mapping it on *Siracusa*. In Section V, we present the *Siracusa* MCU platform. Section VI presents and discusses the end-to-end deployment results, comparing them to the state-of-the-art. Finally, Section VII concludes this article, summarizing the results and contributions.

## II. RELATED WORK

This section gives an overview of the state-of-the-art on EFMs, focusing on the developments toward improvements in energy efficiency and model size and tools to deploy the DNNs on the extreme edge devices.

### A. Small Foundation Models

Recently, the development of the decoder-only large language models (LLMs), such as Llama [1], and Mixtral [2], and their associated ML pipelines led to a new model type: the FM.

EFMs are pretrained LLMs, which can be fine tuned for the downstream tasks at a fraction of the cost of pretraining,

<sup>1</sup>We will open-source all code required to reproduce our experiments under <https://github.com/pulp-platform/deeploy>.

175 making them particularly relevant for the domain specializa-  
 176 tion. However, LLMs often contain several billion parameters,  
 177 requiring GiB of storage space, making them incompatible  
 178 with the extreme edge inference.

179 Addressing this gap, the emerging field of SLMs has gained  
 180 significant traction in the last year. The aim of SLMs is to  
 181 compact LLMs down to tens to hundreds of MiB [8], [9],  
 182 mirroring the evolution of compression of CNNs [17] over the  
 183 past decade.

184 This paradigm shift toward compact FMs is particularly  
 185 interesting for TinyML applications. Incorporating smaller  
 186 FMs, like SLMs, into the embedded devices may enable a new  
 187 wave of intelligent, responsive, and autonomous devices built  
 188 on the EFM. Such systems could bridge the gap between  
 189 the human-understandable inputs, such as text and performing  
 190 high-level planning and low-level control tasks [18] and make  
 191 such advanced capabilities available at the edge, embodied in  
 192 robots, appliances, and wearable devices.

193 In this work, we contribute to the growing field of SLM  
 194 and EFM research and aim to lay the foundation for truly  
 195 embedded SLMs by providing a foundational deployment flow  
 196 that supports a wide range of FMs, from the autoregressive  
 197 decoder-only ones to the encoder-only ones.

### 198 B. Neural Network Deployment for Extreme Edge Devices

199 Building on the trends of the model quantization and  
 200 compression, as well as research into more computation-  
 201 ally efficient DNNs [25], DNN inference on the mobile  
 202 and embedded devices has become a flourishing field of  
 203 research [13], [14], [26]. While model deployment on the  
 204 mobile devices like smartphones follows similar approaches to  
 205 the server-scale deployment, relying on the ample compute and  
 206 memory resources, hardware-managed caches, and operating  
 207 systems to carry out task scheduling available to this class  
 208 of devices, deeply embedded devices face much more severe  
 209 constraints in deployment. This is especially true for the  
 210 new generation of MCU-class devices focusing on the AI  
 211 applications. In contrast to their predecessors, these MCUs  
 212 feature multicore compute clusters, DNN accelerators, and  
 213 on-chip memory of up to 10 MiB, split into multiple software-  
 214 managed memory hierarchy levels [13], [14], [27].

215 To optimally leverage the compute capabilities of such  
 216 complex systems, network deployment must simultaneously  
 217 optimize the execution schedule and tiling of operators and  
 218 orchestrate overlapping memory transfers using DMAs to  
 219 achieve low data marshaling overheads and high compute  
 220 utilization. While modern top-down compilers like MLIR and  
 221 TVM [15], [16] allow integration of most common instruction  
 222 set architectures (ISAs) and accelerator APIs, their focus is  
 223 not on meeting the stringent memory constraints of this class  
 224 of TinyML devices. Prior work like Dory [28], CoSa [29],  
 225 and others have addressed these challenges for CNNs by  
 226 focusing on the operator tiling to fit the target’s memory  
 227 constraints. However, these approaches assume a single-  
 228 cluster memory hierarchy, with undivided memory at each  
 229 level, and a simple lifetime model for the network tensors,  
 230 which are fundamentally stateless across the inference rounds.

231 These simplifying assumptions do not hold for the complex  
 232 heterogeneous multiaccelerator hardware and advanced SLM  
 233 networks [30], [31].

234 Moving beyond these prior works, we propose a novel  
 235 constraint programming algorithm that enables co-optimizing  
 236 tiling and memory allocation, which overcomes the limitations  
 237 of the previous approaches by supporting the data flows with  
 238 complex lifetimes (e.g., KV caching) as required by EFMs.

## 239 III. DEEPLOY

240 In this section, we provide an overview of the Deeploy  
 241 compilation flow. In contrast to most state-of-the-art compil-  
 242 ers for DNNs, which lower DNN representations top-down  
 243 into predefined primitives that need to be implemented by  
 244 each backend [15], [16], [32], Deeploy employs a bottom-  
 245 up compilation approach, where the compiler implements  
 246 networks by composing the user provided C kernels, extending  
 247 them with the code generation passes to implement tiling and  
 248 memory allocation. This bottom-up approach to compilation  
 249 provides three key advantages: first, it supports reusing hand-  
 250 optimized kernel libraries commonly available for most ISAs  
 251 and accelerators. Second, it can be easily extended to support  
 252 highly customized nonstandard compute platforms, including  
 253 heterogeneous SoCs featuring multiple accelerators for which  
 254 a low-level compiler backend may not exist. Third, it allows  
 255 easy integration of novel operators found in emerging the  
 256 Transformer architectures without invasive modifications to the  
 257 deployment flow.

258 Deploy is organized in three building blocks; the *front-*  
 259 *end* validates and transforms the graph representation into  
 260 a representation that suits the platform and assigns kernel  
 261 templates to each operator. The *midend* performs all the  
 262 tiling and static memory allocation computations, guaranteeing  
 263 that the computed program schedule may execute without  
 264 unscheduled runtime memory spills. Finally, the *backend* uses  
 265 the optimized graph representation generated in the *frontend*,  
 266 and the generated tiling schedule and memory allocation  
 267 map generated in the *midend* to create the executable code  
 268 through a series of code generation passes. All the deployment  
 269 targets share the same execution flow, and Deeploy uses a  
 270 configurable platform abstraction, the *backend*, which allows  
 271 it to steer the operators’ mapping, optimization, and lowering  
 272 according to the platform’s configuration. An overview of the  
 273 Deeploy execution flow is shown in Fig. 1.

### 274 A. Data Structures

275 Deeploy distinguishes between three types of buffers:  
 276 1) *variable buffers*; 2) *transient buffers*; and 3) *constant*  
 277 *buffers*. *Variable buffers* represent tensors that contain data  
 278 that is not constant at compile time, i.e., network inputs,  
 279 outputs, and intermediate activations. *Constant buffers* rep-  
 280 resent compile-time constant data used in inference, i.e.,  
 281 network weights and other network parameters. Finally, *tran-*  
 282 *sient buffers* represent scratchpad memory locations for the  
 283 kernel execution, e.g., *im2col* buffers for the convolution  
 284 kernels [33], [34], or reorder buffers for efficient transposition  
 285 kernels. Typically, the amount of space used in *transient*

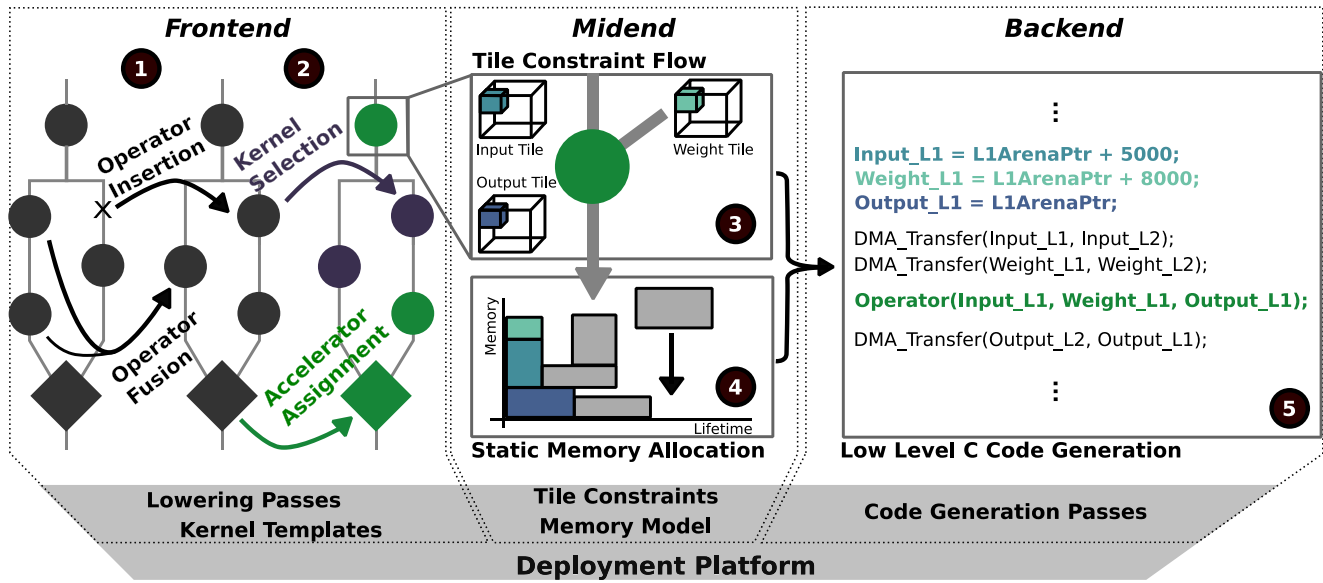


Fig. 1. Overview of the Deeploy execution flow. Steps ① and ② are part of the *frontend*. In the first step, the graph is modified by fusing and inserting platform-specific operators, for example, transposition operators, to match the data layout requirements. In the second step, the datatypes for every tensor are inferred, the accelerator target is chosen, and the kernel templates are selected. The first step in the *midend*, step ③, is the TCF, which computes geometrical constraints for the tile sizes of each tensor, adding them to a CP. The resulting tensor size variables are translated into a 2-D bin packing problem in step ④. The solution of the co-constrained tiling and static memory allocation problem is computed by the ORTools CP-SAT solver and finally processed in step ⑤ in the *backend*. Step ⑤ generates platform-specific C code exploiting DMA transfers. Each step of the execution flow is highly configurable through the *Backend* object.

286 *buffers* depends on the operator’s parametrization, distinguishing 318  
 287 them from the *variable buffers*. In contrast to simpler DNN 319  
 288 topologies, EFMs employ data structures that require advanced 320  
 289 allocation strategies, such as the *KV* caches of the autoregressive 321  
 290 SLMs, as they have more complex buffer lifetime 322  
 291 requirements than the intermediate tensors found in CNNs. 323  
 292 Addressing these constraints requires a more sophisticated 324  
 293 management of the buffers’ lifetime and memory allocation 325  
 294 than in the other deployment tools targeting the extreme edge 326  
 295 devices [28], [29]. 327

296 The distinction between the global and local section buffers 328  
 297 is relevant for the code generation; the global objects are 329  
 298 allocated as the global C variables, while the local objects 330  
 299 are only accessible in the inference code. As such, the global 331  
 300 variables are alive throughout an inference execution, while the 332  
 301 local variables are allocated and deallocated as the network’s 333  
 302 execution schedule requires. 334

### 303 B. Frontend

304 Deeploy’s *frontend* is designed around ingesting quantized 335  
 305 open neural network exchange (ONNX) graphs produced by 336  
 306 DNN and Transformer quantization tools like Quantlib [35]. 337  
 307 Deeploy implements a configurable lowering pass system 338  
 308 based on the pattern matching of the ONNX graphs to enable 339  
 309 efficient and customizable graph-lowering strategies. Each 340  
 310 lowering pass consists of an user-defined replacement function 341  
 311 and a *source pattern*, which describes the subgraph that should 342  
 312 be replaced. Using the replacement function, each lowering 343  
 313 pass uses the matched subgraph to generate a *target pattern*, 344  
 314 which replaces the *source pattern*. Using this system, the 345  
 315 first processing step in the *frontend* is transforming the input 346  
 316 graph into a custom, platform-specific ONNX dialect using 347  
 317 lowering passes provided by the *backend*. The user further 348  
 349  
 350

318 defines operator mappings between the custom operators 319  
 319 and the engines available in the target platform to control 320  
 320 the code generation on the level of individual operators. 321  
 321 Common TinyML kernel libraries like CMSIS-NN and PULP- 322  
 322 NN [33], [34] offer kernels for the fused linear operators 323  
 323 and activations, which can be lowered into by matching 324  
 324 pairs of linear operators and quantization operators. Besides 325  
 325 operator fusion optimization passes, Deeploy also supports the 326  
 326 minimization and insertion of the data marshaling operators 327  
 327 like transpositions to match the data layout requirements of 328  
 328 the kernel libraries. An example of such an operator insertion 329  
 329 pass is adding transpositions operators to optimize the data 330  
 330 layout of the *B* matrix for the general matrix multiplication 331  
 331 (GEMM) kernels of type  $Y = \alpha AB + \beta C$  for better data access 332  
 332 locality. 333

333 The second step after transforming the input graph into 334  
 334 the platform-specific dialect in the *frontend* is parsing, during 335  
 335 which every operator in the network is analysed to construct 336  
 336 an initial context of buffers used in the network’s execution, 337  
 337 and *type inference and kernel selection* where every buffer 338  
 338 in the context is assigned a type. The types used in Deeploy 339  
 339 correspond to the standard C types (e.g., *int8\_t*, *float32*) 340  
 340 or custom data types, depending on the kernels used by 341  
 341 the *backend*. To guarantee a valid type assignment, Deeploy 342  
 342 propagates type information top-to-bottom. The user must 343  
 343 only provide the input types for every graph’s input tensor to 344  
 344 achieve this; then, using this information, Deeploy matches the 345  
 345 input types of each operator with one of the kernel signatures 346  
 346 provided by the *backend*. 347

347 The final result of the *frontend* is an assignment of low-level 348  
 348 kernel templates to every operator in the lowered platform- 349  
 349 specific ONNX, which satisfies the type constraints imposed 350  
 350 by the network’s operators. 351

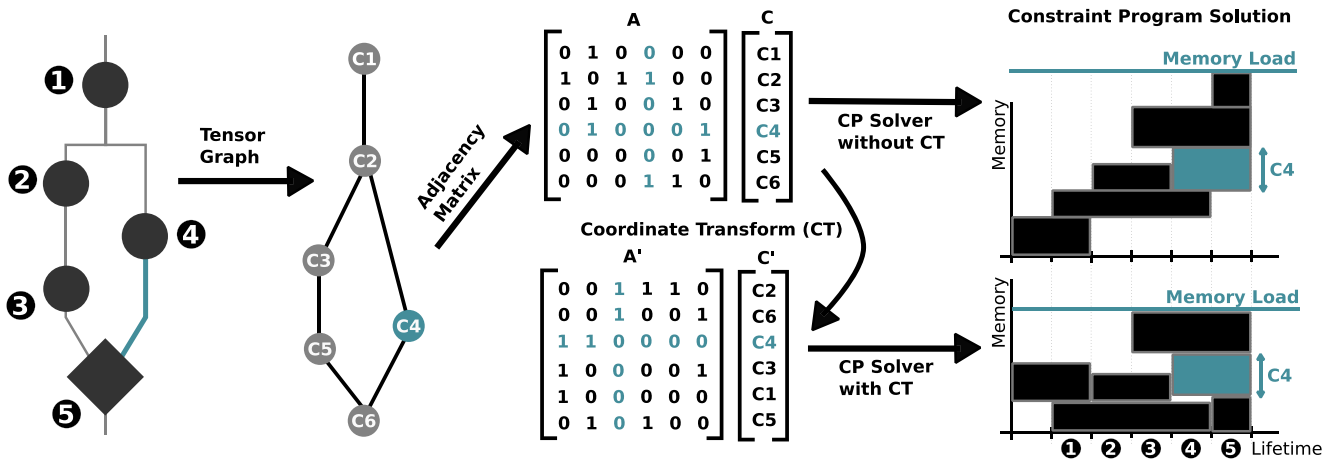


Fig. 2. Example of the co-optimization of tiling and static memory allocation algorithm for one memory level in Deeploy. First, the lifetime of each tensor in the graph is calculated under the execution schedule shown on the left. Next, the memory scheduler constructs an adjacency matrix of the tensor graph and extracts the cost vector from the TCF shown in the middle. Finally, Deeploy applies a coordinate transform within the CP. On the right-hand side, the 2-D bin packing solution is presented with the naive solution on the top, and the solution found by Deeploy is shown below.

### 351 C. Midend

352 The second stage of Deeploy’s execution flow, the *midend*,  
 353 receives the platform-specific ONNX graph and the kernel  
 354 assignment for each operator from the *frontend*. The *midend*’s  
 355 purpose is to perform all the optimization operations required  
 356 to generate the low-level optimized C code for the target  
 357 platform in the *backend*. The *midend* is divided into two  
 358 optimization steps: 1) *memory level annotation* and 2) *tiling*  
 359 *and memory scheduling*. To model the CP used to compute  
 360 the tiling and static memory allocation solution, Deeploy uses  
 361 Google’s ORTools.

362 1) *Memory Level Annotation*: The memory level annota-  
 363 tion step annotates every buffer in the compilation context  
 364 with a memory hierarchy level. The motivation for defining  
 365 the storage location of every tensor is to model the code  
 366 generation constraints closely to the hardware; most embedded  
 367 systems designed for the TinyML applications use multiple  
 368 memory or cache levels [13], [14] to optimize the tradeoff  
 369 between the storage density and the memory access latency.  
 370 While Deeploy supports the tiling of buffers, directly assign-  
 371 ing buffers’ memory levels to lower cache levels can lead  
 372 to performance improvements. When targeting accelerators  
 373 that would otherwise be limited by the available bandwidth  
 374 toward higher-level caches, controlling memory allocation has  
 375 a significant performance impact [13].

376 2) *Tiling*: The second processing step in the *midend* is  
 377 *tiling and memory scheduling*. For every kernel template  
 378 chosen in the *frontend*, the target platform must specify a tile  
 379 constraint (TC). The TC models the geometric and platform-  
 380 specific constraints for tiling an operator. For a tiling solution  
 381 to be correct, all the geometric constraints must hold. For  
 382 example, the spatial dimensions of a Softmax activation’s  
 383 output tile must be the same as its input tile’s dimensions.  
 384 As such, the geometric constraints do not depend on the  
 385 implementation of an operator. While it is possible to tile  
 386 the large tensor operators down to single instructions when  
 387 targeting the processor cores, the same does not hold for the

accelerators. Specifying TCs and platform-specific constraints 388  
 on a per-kernel basis is especially important for handling the 389  
 tiling problem for the loosely coupled accelerators since they 390  
 typically only support specific dimensions to be tiled, owing 391  
 to their specialized datapaths [13], [14]. 392

393 Similarly to the *type inference and kernel selection* flow,  
 394 the tile constraint flow (TCF) is applied top-to-bottom through  
 395 the execution schedule of the network, adding the geo-  
 396 metric and platform-specific TCs of every operator to the  
 397 CP. Furthermore, the TCF adds one symbolic variable per  
 398 dimension per tensor in the network to the CP and a sym-  
 399 bolic variable for every tensor, representing its size as the  
 400 product of all the dimension variables. Using this formulation,  
 401 the solution of the CP represents the size of the largest tile.  
 402

403 3) *Memory Scheduling*: After the geometrical constraints  
 404 of every mapped kernel template in the network are collected  
 405 and added to the CP, the Deeploy’s memory scheduler calcu-  
 406 lates the lifetime of every tensor in the network over the user  
 407 provided execution schedule of the ONNX graph as shown in  
 408 Fig. 2. As previously mentioned, this is an essential step for  
 409 the autoregressive Transformers that must accommodate short-  
 410 lived tensors (e.g., intermediate activations and residuals) and  
 411 long-lived buffers (such as KV caches).

412 Deeploy’s memory scheduler computes a tiling path using  
 413 the *backend*’s memory hierarchy model to assign a sequence  
 414 of memory transfers through the different memory levels.  
 415 Using the calculated lifetimes and the tensor’s size variable  
 416 computed before, the memory scheduler models the problem  
 417 of computing a static memory allocation schedule as a 2-D  
 418 bin packing problem [30], [36], where the horizontal axis  
 419 represents the lifetime, and the vertical axis represents the  
 420 memory address space.

421 Similar to the other state-of-the-art algorithms [30],  
 422 Deeploy’s scheduling CP works with Tetris scheduling intro-  
 423 duced in Tetrisched [37], where the memory buffers are  
 424 scheduled one after another, adding to the maximum load  
 425 of each of their lifetime’s bins. To solve the tiling and

426 allocation problem in a single shot, the memory allocation of  
 427 each buffer is coupled to the tiling solution, which requires  
 428 expressing the order in which they are scheduled within the CP  
 429 as well.

430 The first step to modeling the memory allocation problem is  
 431 to pick a random schedule of the memory buffers and compute  
 432 the adjacency matrix  $A$  of the tensor graph. We collect the  
 433 memory size of each buffer, represented as an integer variable  
 434 of the CP, in a cost vector  $C$ . For any permutation matrix  $P$ ,  
 435  $A' = P \times A \times P^T$  is a valid adjacency matrix with associated  
 436 cost vector  $C' = P \times C$ . A valid  $N \times N$  permutation matrix  
 437 can be expressed as

$$\begin{aligned}
 438 \quad & p_{i,j} \in [0, 1] \quad \forall i, j \in [0, N - 1] \\
 439 \quad & \sum_{i=0}^{N-1} p_{i,j} = 1 \quad \forall j \in [0, N - 1] \\
 440 \quad & \sum_{i=0}^{N-1} p_{j,i} = 1 \quad \forall j \in [0, N - 1].
 \end{aligned}$$

441 Next, the total memory load is computed iteratively using  
 442  $A'$  &  $C'$ : since we use Tetris scheduling, we add each buffer's  
 443 memory size to the size of the final scheduled buffer whose  
 444 lifetime overlaps. We use a vector of intermediate variables  
 445 containing one entry for each buffer,  $H$ , representing the  
 446 memory load in the lifetime region of each buffer. The vector  
 447  $H$  is computed as follows:

$$\begin{aligned}
 448 \quad & H_0 = 0 \\
 449 \quad & H_j = \max_{i=0 \dots j-1} (A'[j, i] \cdot H_i) + C'_j.
 \end{aligned}$$

450 The total worst-case memory load for all the execution steps  
 451 is then computed as memory load =  $\max_{i=0 \dots N} (H_i)$ .

452 In contrast to the other static memory schedule algorithms,  
 453 which focus on calculating an optimal solution for the memory  
 454 blocks of the fixed size, our algorithm combines the constraints  
 455 on the tile sizes and memory layout calculation into a single  
 456 CP; this allows Deeploy to simultaneously optimize static  
 457 memory allocation as well as tile sizing to control memory  
 458 use during the entire inference process, which is critical to  
 459 matching the memory constraints of extreme-edge SoCs with  
 460 the complex buffer lifetime requirements of Transformers.  
 461 An overview of the co-constrained tiling and static memory  
 462 allocation algorithm is shown in Fig. 2.

#### 463 D. Backend

464 Every kernel template picked in the *frontend* is assigned  
 465 a list of code generation passes by the *backend*. Each code  
 466 generation pass operates on a code segment, starting from  
 467 the original kernel template, and may add to or modify its  
 468 code segment. Besides enabling integration of custom passes,  
 469 Deeploy offers the standard code generation passes required  
 470 for generating the correct code, e.g., memory allocation and  
 471 deallocation generation, which inserts calls to heap-based  
 472 allocators or sets pointers to predefined memory locations  
 473 calculated during *tiling and memory scheduling*.

474 An essential set of code generation passes is centered around  
 475 generating closures for the code segments. In the context of

### Kernel Signature

```

// Function signature
void gemv_s8_s8(int8_t* input, int8_t* weight,
                int32_t* bias, int8_t* output,
                uint16_t M, uint16_t N, uint16_t O);

```

### Kernel Template

```

// Kernel Template
gemv_s8_s8({A}, {B}, NULL, {C}, 1, {N}, {O});

```



### Cluster Offloading

#### Global Definitions

```

typedef struct {
    int8_t* A;
    int8_t* B;
    int8_t* C;
} GEMV_closure_args_t;
void GEMV_closure(void* GEMV_closure_args){
    GEMV_closure_args_t* args =
        (GEMV_closure_args_t*) GEMV_closure_args;
    int8_t* _A = args->A; int8_t* _B = args->B;
    int8_t* _C = args->C;

    gemv_s8_s8(_A, _B, NULL, _C, 1, {N}, {O}); }

```

#### Kernel Replacement

```

// GEMV Closure Call
GEMV_closure_args_t GEMV_closure_args = {
    .A = {A}; .B = {B}; .C = {C};
};
// Parallelize Closure over eight cores
pi_cl_team_fork(GEMV_closure, &GEMV_closure_args, 8);

```

Fig. 3. Bottom-up offloading closure generation for a GEMV kernel. All the arguments that refer to the nonglobal *variable buffer's* or *constant buffer's* are captured and used to generate a closure struct typedef and a closure function that unpacks the argument struct and calls the original kernel. Finally, the kernel template is replaced with a function *pi\_cl\_team\_fork*, which takes the newly generated closure as an argument and offloads its execution to all the eight cluster cores.

Deploy, closure generation consists of three parts: 1) the  
 476 closure function itself, which encapsulates a code segment; 2)  
 477 the closure environment, which contains every free variable  
 478 used within the code segment and must be passed to the  
 479 closure function; and 3) the closure invocation, which is either  
 480 an offloading function or a call to the closure function.  
 481

Deploy implements closures as the standard *C* functions by  
 482 generating a function call around the target code segment and  
 483 passing the closure environment as a struct pointer. Deeploy  
 484 captures the relevant free variable expressions by analysing  
 485 the abstract syntax tree (AST) of the underlying code segment  
 486 using the Mako templating library [38]; since the function  
 487 signature of the kernel template is known to Deeploy, it can  
 488 extract arguments used in the kernel template that refer to local  
 489 buffers, and pass them to the closure using an argument struct.  
 490 During code generation, the closure generation pass hoists the  
 491 closure function definition into the global context, inserts code  
 492 for constructing the argument struct and returns the function  
 493 call to the hoisted closure as the new code segment for the  
 494 subsequent code generation passes.  
 495

An important application for the Deeploy's closures is to  
 496 facilitate operator offloading, which is required for program-  
 497 ming processor-based accelerators like compute clusters or  
 498

loosely coupled, memory-mapped accelerators like NPUs. An example of closure generation for the operator offloading to the octa-core cluster is shown in Fig. 3.

Tiling code generation is implemented as a pass as well. Deeploy supports DMA engines and uses them in tiling code generation to move tiles between different memory hierarchy levels according to the tiling solution computed in *tiling and memory scheduling*. To hide the latency of DMA transfers, Deeploy can configure tiling for the operators to use double-buffering, which constrains the tiling solution to reserve twice the required space for every input and output tile. During code generation, Deeploy schedules the data fetching and writeback to occur in parallel with the kernel execution to minimize the latency.

#### IV. TINYSTORIES LLAMA MODEL

As a concrete example of our deployment flow for the next-generation EFMs, we quantize and deploy an SLM on a heterogeneous MCU, Siracusa, introduced in Section V. We chose a Llama2 model pretrained on the TinyStories dataset [8] from HuggingFace,<sup>2</sup> with a hidden size  $d_m = 64$ ,  $h = 16$  parallel attention heads,  $N = 8$  layers, and an intermediate size  $d_{ff} = 256$  for the feed-forward layer. The model architecture is shown in Fig. 4. Note that, however, any SLM fitting the memory constraints of the target platform can be deployed with the same flow.

Like all the other decoder-based language models, the Llama model we use in this work has two fundamental inference modes, which we refer to as the autoregressive and parallel inference modes, and generates its response in two distinct phases, the *prompting* and *generation phases*; the prompting phase ingests the initial sequence of the user input tokens, whereas the generation phase generates the model’s output tokens autoregressively.

##### A. Prompting Phase

Inferences follow a two-pass regime. First, the text input is translated into a sequence of tokens, typically referred to as the prompt. The prompt can have an arbitrary sequence length  $S_p$ , up to the size of the context window of the model.

In the first pass of the model, the prompt is processed to produce the first output token. Since, all the tokens of the prompt are available ab initio, the decoder can process them in a parallel single-shot fashion by applying causal masking of the attention matrix [39]. This first pass generates the first token output and the  $K$  and  $V$  matrices, which may be reused in the subsequent *generation phase*. This process parallels the function of encoder layers used in the first Transformer models [39].

##### B. Generation Phase

In the generation phase of the inference process, output tokens are generated one at a time using the previous token outputs as the model’s input. While every step of the generation phase may use the same parallel inference mode described

in the previous Section, doing so would require recomputing all the previous tokens’  $K$  and  $V$  submatrices. Therefore, the  $K$  and  $V$  matrices of the previous inference steps are typically cached in memory to avoid the quadratic cost of recomputing them [39].

As the parallel and autoregressive inference modes require different tradeoffs in memory allocation for  $KV$  caching and storage of intermediate results we deploy them using separate ONNX models which reflect these tradeoffs. For the parallel inference mode we export an ONNX model with a single input and output for the token sequence and outputs for the computed  $KV$  submatrices which are stored for the next generation phase. For the autoregressive inference mode, we use an ONNX model that additionally requires the cached  $KV$  submatrices. While computing outputs using  $KV$  caches is significantly more efficient regarding the absolute number of operations, loading and storing the  $KV$  caches induces significant data movement, and the smaller operator dimensions make the generation phase much more challenging to accelerate.

#### V. DEPLOYMENT PLATFORM

This Section introduces the hardware platform used in this work as a deployment target to deploy the SLM introduced in Section IV and goes over the NPU-specific *backend* implementation in Deeploy.

##### A. Siracusa

Siracusa [13], is a low-power, heterogeneous RISC-V MCU implemented in TSMC 16 nm technology, which is the multi-accelerator SoC targetted in this work. Siracusa is designed for efficient AI inference, which can leverage its dedicated NPU, *N-Eureka*, and generalistic digital signal processing (DSP) tasks, which can exploit both dedicated XpulpNN ISA extensions [33] enabling single instruction multiple data (SIMD) processing of low-precision integers, as well as an accelerator cluster of eight RISC-V cores which enable single program multiple data (SPMD) processing.

To enable single-latency access from the cluster cores to the L1 tightly coupled data memory (TCDM), all the cores and the 16 L1 memory banks are connected through a TCDM interconnect using one 32-bit port each, granting a total memory bandwidth of 256 bit per cycle to the compute cluster. The cluster’s TCDM memory banks are also accessible from the *N-Eureka* accelerator using nine-bank wide, 288 bit accesses. To manage contention on accesses to the single-ported memory banks, Siracusa integrates a lightweight, programmable access arbiter, which allows the set the maximum number of stall cycles for the accelerator; if accesses from the core-side interconnect cause accelerator access to stall for the programmed number of cycles, the arbiter will stall core accesses and grant it to *N-Eureka*.

The *N-Eureka* accelerator uses a mixed-weight-precision bit-serial datapath, which is optimized for executing dense  $3 \times 3$ , depthwise  $3 \times 3$ , and dense  $1 \times 1$  convolution operations with 8 bit activations and 2 bit to 8 bit convolution weights [13]. To support the bit-serial nature of the datapath,

<sup>2</sup><https://huggingface.co/Maykeye/TinyLLama-v0>

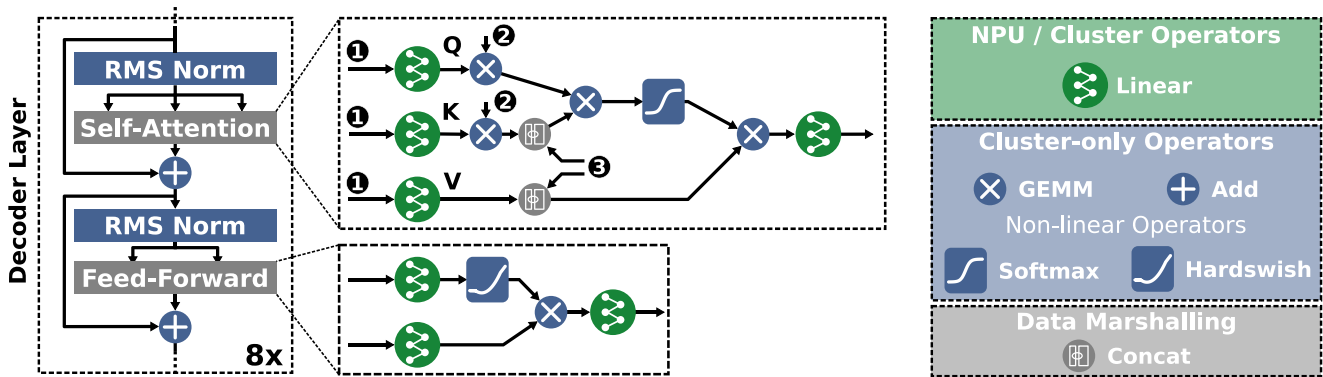


Fig. 4. Overview of the Llama model deployed in this work. The eight decoder layers of the model are shown on the left and consist of an *RMSNorm* - *Self-Attention* - *RMSNorm* - *Feed-Forward* layer stack. Input ① in the self-attention inset corresponds to the token input. Input ② corresponds to the rotational embedding used in Llama models. Input ③ are the *KV* cache inputs used during the autoregressive inference. Notably, during the autoregressive inference, the new row of the *K* and *V* matrices computed on the input token are appended to the *KV* cache.

606 *N-Eureka* requires its weights to be stored in a nonstandard bit-  
 607 interleaved data format, which requires offline transposition,  
 608 padding, and bit shuffling of CNN weight tensors. *N-Eureka*  
 609 is designed as an output-stationary accelerator, opting to cache  
 610 small input tiles and streaming weights. To execute operations  
 611 larger than its internal buffers, it integrates a hardware tiler  
 612 with a programmable number of tiles and strides between  
 613 the dimensions and fixed tile sizes that match the buffer  
 614 sizes. To increase the available memory bandwidth for the  
 615 *N-Eureka*'s weights and minimize off-chip access to fetch  
 616 weights, the cluster integrates a neural memory subsystem  
 617 (NMS), which contains two dedicated 4 MiB memory subsys-  
 618 tems, implemented in static random access memory (SRAM)  
 619 and MRAM technology, respectively, which are designed to  
 620 hold the weights for the *N-Eureka* accelerator and are attached  
 621 through a dedicated 256 bit per cycle weight data port.

622 The compute cluster and *N-Eureka* are located in a shared  
 623 clock domain, the heterogeneous cluster, which communicates  
 624 with the remainder of the SoC, mainly consisting of a  
 625 controller core, 2 MiB L2 memory, and peripherals, through a  
 626 64 bit wide advanced extensible interface bus (AXI) bus, which  
 627 can be used by a DMA integrated within the cluster, to transfer  
 628 the data between the L1 and L2 memories autonomously.

629 While Siracusa is equipped with significant computing  
 630 capabilities through two dedicated accelerators and sizeable  
 631 on-chip memory, deploying an advanced neural network on  
 632 this device is a challenging problem. While weight storage for  
 633 the layers that can be executed on *N-Eureka* is plentiful, all  
 634 the other layers' activation, weight, and output tensors must  
 635 be tiled to fit within 256 KiB of L1 memory. Furthermore,  
 636 memory transfers between L2 and L1 should be orchestrated  
 637 using the DMA to minimize stalling.

### 638 B. Deploy Integration

639 We address the deployment challenges posed by Siracusa's  
 640 heterogeneity through an augmented *Backend* model. This  
 641 section gives an overview of the additions implemented to use  
 642 Deploy for deploying SLMs on Siracusa and, more generally,  
 643 of the modifications needed to support a generic new platform  
 644 in our deployment tool.

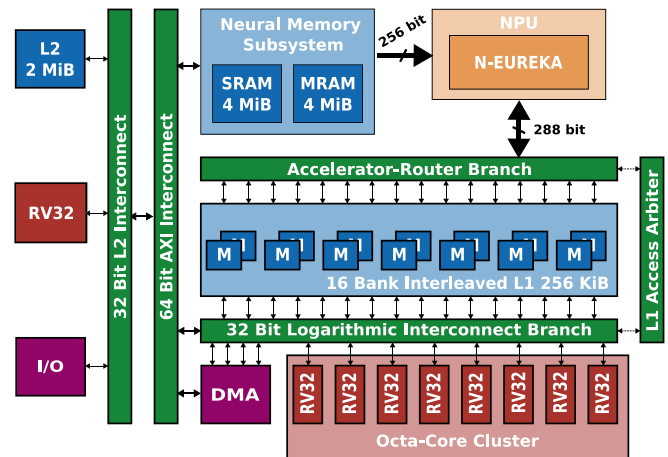


Fig. 5. Overview of the Siracusa SoC featuring its DSP-enhanced octa-core RISC-V cluster and host controller (red), NPU (orange), complex memory hierarchy with two levels of scratchpad memory and a NMS (blue), two arbitrated interconnects toward the L1 memory and an AXI interconnect (green), and peripherals, such as the cluster DMA and chip-level I/O (purple).

645 As Deploy's core primitives are optimized kernels, we  
 646 chose the PULP-NN [33] kernel library, which integrates  
 647 parallel kernels as well as single-core implementations, as  
 648 our target for utilizing the octa-core cluster. The PULP-NN  
 649 kernels focus on efficient implementations of fused linear and  
 650 quantization layers. We support fused layers through lowering  
 651 passes that match the supported operator combinations and  
 652 merge them in the *frontend* of Deploy. We further added  
 653 the fused linear operator TCs, which add the kernel-specific  
 654 constraints besides providing general geometric constraints.

655 We implement function offloading to both the NPU and the  
 656 octa-core compute cluster in Siracusa using the closure system  
 657 as detailed in Section III-D.

658 The *N-Eureka* accelerator provides greater compute capa-  
 659 bilities than the octa-core cluster for the CNN operators,  
 660 achieving a peak throughput in the range of hundreds of  
 661 GOP/s for pointwise and  $3 \times 3$  convolutions. Even though  
 662 SLMs do not employ these types of operations, we add a  
 663 custom *linear layer to pointwise convolution* lowering pass  
 664 that converts GEMM operators with compile-time constant



weight matrices into pointwise convolutions. This method allows us to deploy all the linear layers in Transformer models as shown in Fig. 5, on the NPU.

### C. Deployment Setup

As explained in Section IV, the dual inference modes of the decoder-only models require different deployment strategies, as the autoregressive inference mode requires significant memory for *KV* caching. We deploy two model prototypes to accommodate this difference, one for the autoregressive inference mode and one for the parallel inference mode.

The autoregressive inference mode model uses additional network inputs corresponding to the previous sequences' *KV* caches. Other than that, the deployment setup between both models is equal. We allocate all the graph inputs and outputs as the global *variable buffers* in Siracusa's L2 memory, and annotate all the local *variable buffers* modeling intermediate tensors in L2 as well. In deployment scenarios that use Siracusa's NMS, we allocate all linear layer weights in the NMS but use L2 for all the activations.

Unless stated differently, all the network operators are executed on the cluster and use Deeploy's TCF to generate tiled inference code, which orchestrates transfers of input, weight, and output tensors between the L2 and L1 memories. For the operators executed on the NPU, weights are stored in the NMS in their entirety and ingested by the accelerator without moving them into L1 first, leveraging the increased available bandwidth from the NMS.

## VI. RESULTS

This section discusses the measurement results of deploying the TinyStories SLM on Siracusa and benchmarking results of general Transformer layers. First, we discuss the setup used to measure the performance results on Siracusa. Finally, we present our benchmarking and end-to-end silicon measurements, as well as the profiling experiments of our compiler.

### A. Deployment Evaluation Setup

To evaluate the model's performance in the autoregressive mode and for causally masked parallel inference, we measure each inference step individually with code generated by Deeploy. We start from the empty *KV* caches for causally masked parallel inference and process  $N$  input tokens simultaneously. We start from the *KV* caches of the previous inference step for all the experiments in the autoregressive mode. To calculate the average throughput and energy per token, we take the average over all 256 inference steps.

We report all the power numbers measured on a Siracusa prototype board using a Keysight N6715C DC, supplying all the operating voltages and measuring current. We perform all the experiments under nominal conditions, i.e., 0.8 V supply voltage and 360 MHz operating frequency of the cluster domain. We measure power consumption for every inference by averaging the power consumption of the model run in a continuous loop.

We measure four distinct deployment scenarios. In the first scenario, *single-core deployment*, we only generate code using a single RISC-V core. In the second scenario, *octa-core*

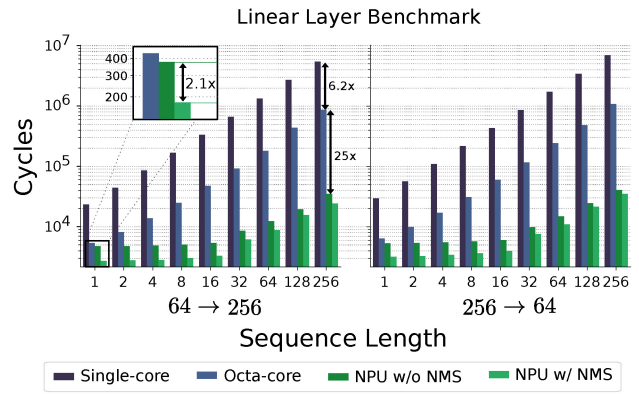


Fig. 6. Performance results for linear layer operators offloaded on *N-Eureka* using Deeploy code generation. The highlighted inset shows that the NMS' added storage and bandwidth leads to performance gains of up to 2.1 $\times$  in memory-bound operator configurations. In large linear layer configurations, the speedup achieved by the NPU is 25 $\times$  compared to the octa-core implementation, and another 1.6 $\times$  when using the NMS for weights.

*deployment*, we generate code using all the eight RISC-V cores of the cluster without using the NPU. In the third scenario, *NPU without NMS deployment*, we generate code using all the eight RISC-V cores and *N-Eureka* without offloading weights to the NMS. In the final scenario, *NPU with NMS deployment*, we generate code using all the eight RISC-V cores and *N-Eureka* with the NMS. We use the Siracusa *Backend* in Deeploy to generate code for all the scenarios.

### B. Microbenchmarking Results

To validate our approach of offloading GEMM operators on *N-Eureka*, we first measure the performance of *N-Eureka* and the RISC-V cluster on the GEMM kernels. Specifically, we study the performance of the Q, K, and V projections in the attention layer and linear layer performance in feed-forward layers for different sequence lengths  $S$  in the parallel inference mode. For the Llama model we study in this article, these projections use dimensions 256  $\rightarrow$  64 and 64  $\rightarrow$  256. Our measurements are shown in Fig. 6.

Transitioning from the single-core to octa-core cluster execution, we measure a performance improvement of 6.2 $\times$ , thanks to the low-overhead parallelization on the cluster cores. Transforming the linear layer operators into pointwise convolutions as explained in Section V-B, enables execution on the NPU, which reduces the latency by 25 $\times$  compared to the octa-core implementation due to the NPU's significant compute resources for the convolution operations. Furthermore, we reduce the data movement by allocating the convolution weights to the NPU's NMS, increasing the effective memory bandwidth available to *N-Eureka*. These optimizations improve the performance, especially on the memory-bound tasks, like the linear layers in attention blocks with low sequence length, by 2.1 $\times$  compared to the NPU execution without the NMS.

We further profile the execution performance of a representative encoder layer as commonly found in nonregressive Transformer models. For our benchmarking, we chose a configuration with the hidden size  $d_m = 64$  and  $h = 16$  parallel attention heads and an intermediate size  $d_{ff} = 256$ , paralleling the decoder layer in Fig. 4. We measure an increase in throughput of 17.8 $\times$  when leveraging the NPU to compute

TABLE I  
CUMULATIVE LATENCY AND ENERGY FOR A 256-STEP INFERENCE OF  
THE SLM ON SIRACUSA USING THE NPU WITH NMS

	Parallel Inference	Autoregressive Inference	Speedup & Energy Reduction Ratio
Cumulative Latency [s]	17.6	0.75	$23 \times$ faster
Cumulative Energy [mJ]	3193	125	$26 \times$ more efficient

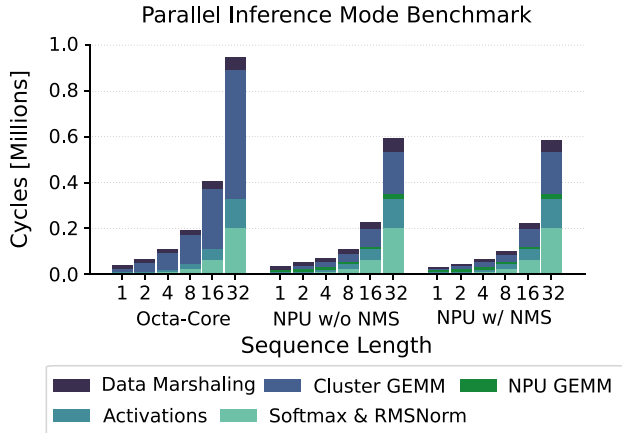


Fig. 7. Cycle breakdown of parallel inference in the studied SLM. Due to the larger contribution of operations from the matrix multiplications using the NMS performance of offloaded GEMM operators increases by 17.8 $\times$ , and end-end-performance improves by 61% for sequence length 32 while maintaining low overheads of only 9%, even when fully leveraging both the cluster and NPU.

759 the linear layers, improving the end-to-end performance for  
760 the encoder layers by 61%. We further quantify the overheads  
761 due to tiling and data marshaling overheads, measuring an  
762 end-to-end overhead of only 9%.

### 763 C. End-to-End Deployment Results

764 We thoroughly evaluate the SLM deployed on Siracusa by  
765 benchmarking the two operating phases required to execute  
766 SLM, namely the *prompting phase* and the *generation phase*.

767 Table I displays the cumulative runtime and energy for exe-  
768 cuting a 256-step inference in the parallel and autoregressive  
769 modes, where KV caching is used. The autoregressive mode  
770 outperforms the parallel mode, achieving a 23 $\times$  speedup and  
771 a 26 $\times$  improvement in energy efficiency. These improvements  
772 directly result from avoiding the costly recomputation of the  
773 KV matrices. Averaging the autoregressive inference mode’s  
774 cumulative latency and energy over 256 steps, we achieve  
775 an average throughput of 340 token per second at an average  
776 energy cost of 490  $\mu$ J/token.

777 Since, the autoregressive mode maximizes the data reuse  
778 across the whole inference process, this mode can be consid-  
779 ered both during the *prompting* and *generation* phases detailed  
780 in Section IV. However, this strategy leads to the suboptimal  
781 results as running in the parallel mode for the *prompting*  
782 phase enables better utilization of the NPU without excessive  
783 recomputation of the KV matrices, as tokens are not fed back  
784 in this phase.

785 The parallel inference mode’s performance for the SLMs  
786 studied in this work follows the trend of the benchmark shown  
787 in Fig. 7. While we benchmark the end-to-end performance of

788 the decoder-only models in this work, the results in Fig. 7 also  
789 apply to the encoder-based transformer models, as the parallel  
790 inference mode is equivalent to the encoder layer execution  
791 in such networks. In the autoregressive mode the speedup  
792 achieved by employing the NPU is only 19%, which can be  
793 attributed to the mode’s smaller operator sizing, leading to  
794 stalling of the accelerator due to the reconfiguration overheads.  
795 Additionally, the average proportion of time spent for the data  
796 marshaling is 40% for the autoregressive versus just 14% for  
797 the parallel modes, underlining the memory access intensity  
798 inherent to KV caching, which drastically reduces the number  
799 of computations leading to reduced arithmetic intensity. A  
800 detailed analysis of runtime and breakdown of the operator  
801 intensity for the end-to-end autoregressive inference is shown  
802 in Fig. 8 plots ① and ②.

### 803 D. Deployment Overheads

804 An important metric for the quality of generated code is  
805 the utilization of the system’s compute engines. To profile  
806 the quality of our code, we measured the overheads incurred  
807 by Deeploy for each autoregressive inference step in *NPU*  
808 *without NMS deployment* and *NPU with NMS deployment*  
809 shown in Fig. 8, plot ③. The main difference between the  
810 two scenarios is whether Siracusa’s NMS is used for the  
811 compile-time constant GEMM weights. While the reduction in  
812 overheads decreases from 33% to 7% with increasing sequence  
813 lengths and arithmetic intensity, the weight memory drastically  
814 reduces the relative time spent on the data movement in the  
815 first steps of inference. This reduction of overheads is a crucial  
816 advantage of the bottom-up compilation approach employed  
817 by Deeploy; while the other compilers might not consider  
818 low-level architectural features like memory hierarchy or only  
819 expose a simplified model, Deeploy allows complete control  
820 over the memory allocation and code generation to leverage  
821 the knowledge of the target architecture fully.

### 822 E. Comparison With TinyML Compilers

823 While we designed Deeploy to deploy the state-of-the-art  
824 and emerging SLMs, we also report the results on more classical  
825 CNN and artificial neural network (ANN) workloads as defined  
826 in the MLPerf tiny benchmark [41]. We compare Deeploy with  
827 the state-of-the-art open-source Dory tool [28] using the same  
828 open-source CNN kernels for PULP MCUs [33] we used in this  
829 work. To ensure a fair, compiler-focused comparison, we do not  
830 use the NMS or the NPUs of Siracusa. In this mode, both the  
831 compilers only deploy cluster kernels with equivalent memory  
832 constraints. As a third data point, we add measurements of  
833 Deeploy-generated code on Siracusa when using the NMS and  
834 NPU. Our results are shown in Table II. We find that Deeploy  
835 generates code with an equivalent latency of Dory up to 1%  
836 of variation, underlining that even though Deeploy chooses a  
837 more general compilation approach than Dory, it does not incur  
838 any performance penalties.

### 839 F. Comparison With the State-of-the-Art

840 Currently, most efforts on EFM deployment target mod-  
841 els with more than a billion parameters on high-end

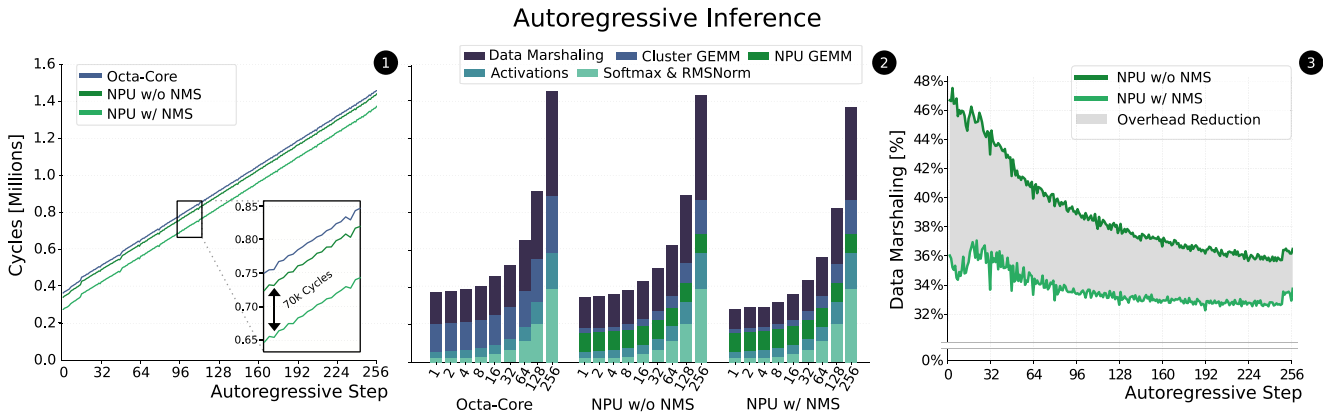


Fig. 8. Performance results of end-to-end autoregressive inference. Plot ① shows the runtime of each autoregressive inference step in three scenarios corresponding to *octa-core deployment*, *NPU without NMS deployment*, and *NPU with NMS deployment*. The plot shows that the autoregressive inference on Siracusa is highly memory-bound in all the scenarios, which is due to the transfer of the KV caches between L2 and L1; *NPU with NMS deployment* reduces the runtime of every step by approximately 70k cycles since the weights are stored untiled in the NMS, reducing the required L2 to L1 data transfers. The second plot ② shows a breakdown of the runtime in the different operators of the network and data marshaling overheads. Evidently, the higher compute throughput of *N-Eureka* is unused due to the overall memory-boundedness. Finally, plot ③ shows that the data movement overhead reduction afforded by the NMS decreases with increasing sequence lengths as the overhead of transferring KV caches increases.

TABLE II  
LATENCY RESULTS OF DORY AND DEEPLOY ON THE MLPERF TINY BENCHMARK, RUNNING ON SIRACUSA AT A CLOCK FREQUENCY OF 360 MHz

Benchmark	Siracusa w/o NPU Dory	Siracusa w/o NPU Deeploy	Siracusa w/ NPU Deeploy
DS-CNN	1.4 ms	1.4 ms	0.39 ms
MobileNetv1	5.6 ms	5.6 ms	0.69 ms
ResNet	3.7 ms	3.7 ms	0.60 ms
ToyAdmos	0.24 ms	0.24 ms	0.11 ms

microprocessors (MPUs) and embedded processors, such as the I.MX95 or NVIDIA Orin or mobile phone chips, featuring multi-GiB external memories and multi-W power envelopes [42], [43]. Even though our performance and efficiency are extremely competitive, quantitative comparisons against these deployments would be unfair in our favor as we target much smaller SLMs.

Considering SLMs in the 100 s million parameters range, we compare our implementation on Siracusa with another small-scale Llama model for the edge devices, MobileLLM, by Liu et al. [44]. Liu et al. deploy a 125 MParameter SLM on an iPhone 13 featuring an A15 Bionic chip in 5-nm technology using the highly optimized metal performance shaders (MPSs) backend for the Apple devices, achieving a throughput of 64 token per second. While their paper does not profile the exact energy consumption of their models during inference, Liu et al. optimistically estimate the energy consumption of their setup with 12.5 mJ per token. Compared to this estimate on the iPhone 13’s A15 processor, the implementation of our SLM on the Siracusa MCU uses 26× less energy per token while achieving 5× more throughput, for a total 130× higher energy efficiency. When normalizing throughput with the number of operations per token of their network, we find that they achieve an equivalent of 4800 TinyStories Llama tokens per second. Under this estimate, our end-to-end energy efficiency on Siracusa implemented in an older 16 nm TSMC technology node is 1.7× higher.

A comparison with a similar-scale (10 s million parameters) model as ours is possible against the *llama2.c* [45] implementation of the TinyStories-15M model on a Samsung Galaxy Watch 4, demonstrated to achieve 22.1 token per second [46] using an Exynos W920 dual-core ARM Cortex-A55 processor [47]. Neglecting the power consumption of dynamic random access memory (DRAM) accesses, only considering a power consumption of 300 mW per core in Samsung 5 nm technology [48], we estimate the power consumption during inference as 600 mW. Under this assumption, the Galaxy Watch 4 achieves an energy efficiency of 27 mJ per token, 55× lower than ours. Normalizing for the operations per token, our energy efficiency is 13.4× greater, even though the Exynos W920 is implemented in an advanced Samsung 5 nm technology node.

## VII. CONCLUSION

In this work, we presented Deeploy, a novel compiler for DNNs allowing broad customizability of deployment flows. We presented the integration of Siracusa, a heterogeneous RISC-V SoC featuring an octa-core compute cluster and an NPU. We demonstrate the deployment of a SLM trained on the TinyStories dataset on Siracusa, achieving a state-of-the-art throughput of 340 token per second at an average energy cost of 490 μJ per token in autoregressive inference mode by efficiently leveraging the on-chip KV caching.

We further analysed the efficiency of our generated code via microbenchmarks, achieving the data marshaling overheads of only 9% on the Transformer encoder layers, even when fully utilizing both the cluster cores and NPU collaboratively.

Finally, we demonstrated that while the data marshaling overheads are significant in the autoregressive inference mode, the energy savings compared to executing the generation phase of SLM in the parallel mode outweigh this drawback, reducing the energy cost per token by 26× while increasing throughput by 23×.

In the future work, we plan to leverage Deeploy's flexibility to support emerging computer architecture innovations, such as multiaccelerator SoCs integrating compute-in memory (CIM) macros.

## REFERENCES

- [1] H. Touvron et al., "Llama 2: Open foundation and fine-tuned chat models," Jul. 2023, *arXiv:2307.09288*,
- [2] A. Q. Jiang et al., "Mixtral of experts," Jan. 2024, *arXiv:2401.04088*.
- [3] G. Team et al., "Gemini: A family of highly capable multimodal models," Dec. 2023, *arXiv:2312.11805*.
- [4] Y. Chen et al., "EEGFormer: Towards transferable and interpretable large-scale EEG foundation model," Jan. 2024, *arXiv:2401.10278*.
- [5] C. Wang et al., "BrainBERT: Self-supervised representation learning for intracranial recordings," in *Proc. 11th Int. Conf. Learn. Represent., 2022*, pp. 1–22.
- [6] M. Ahn et al., "Do as I can, not as I say: Grounding language in robotic affordances," in *Proc. Conf. Robot Learn., 2022*, pp. 1–34.
- [7] D. Driess et al., "PaLM-E: An embodied multimodal language model," in *Proc. 40th Int. Conf. Mach. Learn.*, vol. 202, 2023, pp. 8469–8488.
- [8] R. Eldan and Y. Li, "TinyStories: How small can language models be and still speak coherent english?" May 2023, *arXiv:2305.07759*.
- [9] P. Zhang, G. Zeng, T. Wang, and W. Lu, "TinyLlama: An open-source small language model," Jan. 2024, *arXiv:2401.02385*.
- [10] Y. LeCun, "Deep learning hardware: Past, present, and future," in *Proc. IEEE Int. Solid-State Circuits Conf.*, 2019, pp. 12–19.
- [11] J. Choquette, "NVIDIA Hopper H100 GPU: Scaling performance," *IEEE Micro*, vol. 43, no. 3, pp. 9–17, May 2023.
- [12] N. Jouppi et al., "TPU v4: An optically reconfigurable supercomputer for machine learning with hardware support for embeddings," in *Proc. 50th Annu. Int. Symp. Comput. Archit.*, 2023, pp. 1–14.
- [13] A. S. Prasad et al., "SiracusA: A 16 nm heterogeneous RISC-V SoC for extended reality with at-MRAM neural engine," Dec. 2023, *arXiv:2312.14750*.
- [14] K. Ueyoshi et al., "DIANA: An end-to-end energy-efficient digital and ANALog hybrid neural network SoC," in *Proc. IEEE Int. Solid-State Circuits Conf.*, vol. 65, 2022, pp. 1–3.
- [15] C. Lattner et al., "MLIR: Scaling compiler infrastructure for domain specific computation," in *Proc. IEEE/ACM Int. Symp. Code Gener. Optim. (CGO)*, 2021, pp. 2–14.
- [16] T. Chen et al., "TVM: An automated end-to-end optimizing compiler for deep learning," in *Proc. 13th USENIX Conf. Oper. Syst. Design Implement.*, 2018, pp. 579–594.
- [17] S. Han, H. Mao, and W. J. Dally, "Deep compression: Compressing deep neural network with pruning, trained quantization and Huffman coding," in *Proc. 4th Int. Conf. Learn. Represent.*, 2016, pp. 1–14.
- [18] R. Firoozi et al., "Foundation models in robotics: Applications, challenges, and the future," Dec. 2023, *arXiv:2312.07843*.
- [19] N. Tekin et al., "A review of on-device machine learning for IoT: An energy perspective," *Ad Hoc Netw.*, vol. 153, Feb. 2024, Art. no. 103348.
- [20] A. Gholami et al., "A survey of quantization methods for efficient neural network inference," in *Low-Power Computer Vision*. Boca Raton, FL, USA: Chapman and Hall, 2022.
- [21] S. Kim et al., "I-BERT: Integer-only BERT quantization," in *Proc. 38th Int. Conf. Mach. Learn.*, 2021, pp. 5506–5518.
- [22] S. R. Jain, A. Gural, M. Wu, and C. H. Dick, "Trained quantization thresholds for accurate and efficient fixed-point inference of deep neural networks," in *Proc. Mach. Learn. Syst.*, 2020, pp. 1–17.
- [23] Y. Li et al., "BRECQ: Pushing the limit of post-training quantization by block reconstruction," in *Proc. Int. Conf. Learn. Represent.*, 2020, pp. 1–16.
- [24] G. Xiao, J. Lin, M. Seznec, H. Wu, J. Demouth, and S. Han, "SmoothQuant: Accurate and efficient post-training quantization for large language models," in *Proc. 40th Int. Conf. Mach. Learn.*, 2023, pp. 38087–38099.
- [25] A. Howard et al., "Searching for MobileNetV3," in *Proc. IEEE/CVF Int. Conf. Comput. Vis. (ICCV)*, 2019, pp. 1314–1324.
- [26] J. Lin, W.-M. Chen, H. Cai, C. Gan, and S. Han, "Memory-efficient patch-based inference for tiny deep learning," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 34, 2021, pp. 2346–2358.
- [27] E. Flamand et al., "GAP-8: A RISC-V SoC for AI at the edge of the IoT," in *Proc. IEEE 29th Int. Conf. Appl.-Specif. Syst., Archit. Process. (ASAP)*, 2018, pp. 1–4.
- [28] A. Burrello, A. Garofalo, N. Bruschi, G. Tagliavini, D. Rossi, and F. Conti, "DORY: Automatic end-to-end deployment of real-world DNNs on low-cost IoT MCUs," *IEEE Trans. Comput.*, vol. 70, no. 8, pp. 1253–1268, Aug. 2021.
- [29] Q. Huang et al., "CoSA: Scheduling by constrained optimization forspatial accelerators," in *Proc. 48th Annu. Int. Symp. Comput. Archit.*, 2021, pp. 554–566.
- [30] M. Maas, U. Beaunon, A. Chauhan, and B. Ilbeyi, "TelaMalloc: Efficient on-chip memory allocation for production machine learning accelerators," in *Proc. 28th ACM Int. Conf. Archit. Support Program. Lang. Oper. Syst.*, 2022, pp. 123–137.
- [31] M. D. Moffitt, "MiniMalloc: A lightweight memory allocator for hardware-accelerated machine learning," in *Proc. 28th ACM Int. Conf. Archit. Support Program. Lang. Oper. Syst.*, 2024, pp. 238–252.
- [32] R. David et al., "TensorFlow lite micro: Embedded machine learning for TinyML systems," in *Proc. Mach. Learn. Syst.*, vol. 3, 2021, pp. 800–811.
- [33] A. Garofalo, G. Tagliavini, F. Conti, D. Rossi, and L. Benini, "XpulpNN: Accelerating quantized neural networks on RISC-V processors through ISA extensions," in *Proc. Design, Autom. Test Eur. Conf. Exhib. (DATE)*, 2020, pp. 186–191.
- [34] L. Lai, N. Suda, and V. Chandra, "CMSIS-NN: Efficient neural network kernels for arm cortex-M CPUs," 2018, *arXiv:1801.06601*.
- [35] M. Spallanzani et al., "QuantLab: A modular framework for training and deploying mixed-precision NNs," in *Proc. TinyML Summit*, 2022, p. 1.
- [36] F. Angiolini, L. Benini, and A. Caprara, "An efficient profile-based algorithm for scratchpad memory partitioning," *IEEE Trans. Comput.-Aided Design Integr. Circuits Syst.*, vol. 24, no. 11, pp. 1660–1676, Nov. 2005.
- [37] A. Tumanov et al., "TetriSched: Global rescheduling with adaptive plan-ahead in dynamic heterogeneous clusters," in *Proc. 11th Eur. Conf. Comput. Syst.*, 2016, pp. 1–16.
- [38] "Mako templates for Python." Accessed: Mar. 24, 2024. [Online]. Available: <https://github.com/sqlalchemy/mako/releases>
- [39] A. Vaswani et al., "Attention is all you need," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 30, 2017, pp. 1–15.
- [40] B. Zhang and R. Sennrich, "Root mean square layer normalization," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 32, 2019, pp. 1–14.
- [41] C. Banbury et al., "MLPerf tiny benchmark," Aug. 2021, *arXiv:2106.07597*.
- [42] E. Ruedas, "LLM pipelines: Seamless integration on embedded devices," presented at the Virtual Event, Mar. 2024.
- [43] X. Chu et al., "MobileVLM: A fast, strong and open vision language assistant for mobile devices," Dec. 2023, *arXiv:2312.16886*.
- [44] Z. Liu et al., "MobileLLM: Optimizing sub-billion parameter language models for on-device use cases," Feb. 2024, *arXiv:2402.14905*.
- [45] A. Karpathy, "Llama2.c." Mar. 2024. [Online]. Available: <https://github.com/karpathy/llama2.c>
- [46] Joey (e/l) [@shxf0072]. @Karpathy Llama2.c Running on Galaxy Watch 4. (Dec. 2023). [Online Video]. Available: <https://t.co/sMPCZM3WE4>
- [47] "Samsung galaxy watch4 and watch4 classic teardown." iFixit. Sep. 2021. [Online]. Available: <https://url.zip/4befa14>
- [48] A. Frumusanu "The snapdragon 888 vs the exynos 2100: Cortex-X1 & 5nm—Who does it better?" Feb. 2021. [Online]. Available: <https://www.anandtech.com/show/16463/snapdragon-888-vs-exynos-2100-galaxy-s21-ultra>