



Agostini, N. B., Haris, J., Gibson, P. , Jayaweera, M., Rubin, N., Tumeo, A., Abellán, J. L., Cano Reyes, J. and Kaeli, D. (2024) AXI4MLIR: User-Driven Automatic Host Code Generation for Custom AXI-Based Accelerators. In: International Symposium on Code Generation and Optimization (CGO) 2024, Edinburgh, United Kingdom, 02-06 Mar 2024, pp. 143-157. ISBN 9798350395099 (doi: [10.1109/CGO57630.2024.10444801](https://doi.org/10.1109/CGO57630.2024.10444801))

This is the author version of the work. There may be differences between this version and the published version. You are advised to consult the published version if you wish to cite from it:

<https://doi.org/10.1109/CGO57630.2024.10444801>

<https://eprints.gla.ac.uk/311284/>

Deposited on 11 December 2023

Enlighten – Research publications by members of the University of Glasgow  
<http://eprints.gla.ac.uk>

# AXI4MLIR: User-Driven Automatic Host Code Generation for

1 **Abstract**—This paper addresses the  
2 efficient generation of host driver code for  
3 AXI-based accelerators targeting linear algebra  
4 important workload in various applications such as  
5 learning and scientific computing. We have  
6 focused on automating accelerator driver code  
7 has been paid to the host-accelerator interface  
8 introduces AXI4MLIR, an extension to the MLIR  
9 framework designed to facilitate the generation of  
10 host-accelerator driver code. With novel  
11 transformations, AXI4MLIR empowers the compiler  
12 with new accelerator features (including their instruction  
13 patterns and exploit the host memory hierarchy).  
14 AXI4MLIR’s versatility across different  
15 problems, showcasing significant CPU cycle savings  
16 (up to 56%) and up to a 1.65× speedup in  
17 optimized driver code implementations. This  
18 link will be included for the camera-ready version.  
19 **Index Terms**—MLIR, AXI, Compiler

## I. INTRODUCTION

21 Given the diminishing performance of  
22 today’s general-purpose computing [1], there has been  
23 renewed interest in exploring custom hardware accelerators.  
24 Accelerators can support architecture-level optimizations that  
25 can increase the performance and efficiency of key applica-  
26 tions [2], [3], [4], [5], [6], [7]. One important class of  
27 applications that can benefit from accelerators is tensor algebra  
28 processing, which is widely used in the domains of machine  
29 learning, scientific computing, and data analytics [8], [9], [10].  
30 Tensor operations tend to be computationally intensive and  
31 require high memory bandwidth, making them suitable for  
32 specialized hardware implementations. Automated tools have  
33 been proposed [11], [12], [13], [14] to help explore new  
34 classes of custom domain-specific accelerators targeting tensor  
35 computations, and are currently the best path available to  
36 obtain performance gains in scientific workloads and machine  
37 learning applications.

38 However, designing and fully exploiting custom hardware  
39 accelerators for tensor operations is not a trivial task. When  
40 co-designing these devices, we need to generate efficient archi-  
41 tectures, and we must optimize the communication between  
42 the host CPU and the accelerator. In particular, the host-  
43 accelerator interaction involves several aspects, including data  
44 transfers, synchronization, and the accelerator’s control flow.  
45 These aspects depend on the characteristics of the host CPU  
46 microarchitecture, the host-accelerator interface, the accelera-  
47 tor design, and the application code. Manually rewriting the  
48 host driver code for each accelerator and application scenario  
49 can be very tedious and error-prone. Furthermore, most of

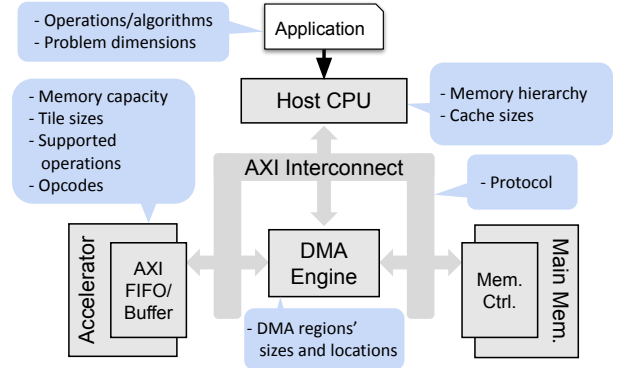


Fig. 1: Typical host-accelerator system design, highlighting (blue color) relevant parameters that should be considered for efficient generation of host-accelerator communication code.

50 the prior work proposing new accelerators [15], [16], [17],  
51 [18], [19] only considers a simple offload model or assumes  
52 that the required data is already placed in the accelerator’s  
53 internal buffers, falling short in providing insights into how  
54 host-to-accelerator transfers should be performed or generated.  
55 Additionally, complex accelerators, exemplified by Google’s  
56 TPUs and Nvidia’s GPUs, benefit from large teams that can  
57 collaboratively engineer dedicated compilers to address some  
58 of these issues. However, smaller development teams may lack  
59 expertise or available time resources to invest in compilers.  
60 Consequently, custom accelerator designers typically imple-  
61 ment driver code and instruction streams manually to validate  
62 and deploy their designs for a subset of synthetic workloads.

63 To implement or generate efficient host-to-accelerator com-  
64 munication, we argue that it is necessary to consider all major  
65 features of a System-on-Chip (SoC). Figure 1 highlights a  
66 typical system using an AXI [20] interconnection between the  
67 CPU and a custom accelerator, which is a common choice  
68 in many designs [21]. To drive the accelerator effectively, the  
69 host-code implementation should exploit features regarding the  
70 CPU, the interconnect, and the accelerator (see Figure 1).

71 To effectively consider each of the key system features  
72 described in Figure 1 while also delivering efficient and  
73 automated CPU-accelerator driver code generation, we  
74 propose **AXI4MLIR**, an extension to the MLIR compiler  
75 framework [22] that enables efficient and automated CPU-  
76 accelerator driver code generation for accelerators targeting  
77 linear algebra applications. AXI4MLIR takes a high-  
78 level application description in the MLIR’s linear algebra  
79 (`linalg`) abstraction [23] as input and introduces custom  
80 MLIR attributes to describe the target accelerator capabilities.  
81 These attributes provide accelerator-specific information to

82 custom transformation passes that can effectively specialize  
 83 and generate accelerator-aware host driver code. Our  
 84 extensions facilitate hardware-software co-design by allowing  
 85 developers to automatically generate driver code with varying  
 86 configurations, more easily explore their design space, and  
 87 use the designed accelerator in applications that can be  
 88 compiled with the MLIR framework. The contributions of  
 89 this work include the following:

- 90 • New MLIR attributes that provide a standardized and  
 91 extensible approach to represent accelerators that can  
 92 implement a range of linear algebra algorithms supported  
 93 by the MLIR *linalg* abstraction.
- 94 • Automated generation of efficient driver code for custom  
 95 accelerators leveraging AXI-based interfaces in host-to-  
 96 accelerator communication.
- 97 • The ability to describe and explore accelerator-specific  
 98 tiling and dataflow strategies for the target linear algebra  
 99 operation, which can improve computation efficiency  
 100 within the accelerator and reduce data movement over-  
 101 heads between the accelerator and CPU.
- 102 • An analysis of our compiler optimizations on a suite  
 103 of benchmarks representing key linear algebra applica-  
 104 tions, demonstrating the effectiveness of our approach  
 105 in achieving significant performance gains (up to  $1.65\times$   
 106 speedup and 56% fewer cache references) when com-  
 107 pared to optimized manual driver code implementations.

108 While leveraging the new attributes of AXI4MLIR, our  
 109 user-directed host code generation is entirely automated by  
 110 the compiler. This provides a significant advantage in terms  
 111 of productivity and maintainability.

## 112 II. BACKGROUND

### 113 A. MLIR

114 MLIR is a compiler infrastructure framework that facilitates  
 115 the creation of domain-specific compilers by providing code  
 116 generators, translators, optimizers, and the infrastructure to  
 117 define subsets of operations that expose well-defined language  
 118 abstractions [22], [24]. Notably, MLIR offers support for  
 119 compilation from various frontends into its infrastructure,  
 120 encompassing frameworks such as TensorFlow, PyTorch, and  
 121 ONNX, as well as languages like Fortran, C, and Mojo.  
 122 In MLIR, a group of operations modeling an abstraction  
 123 is called a *dialect*. Dialects are self-contained intermediate  
 124 representations (IRs) and follow the language rules of MLIR’s  
 125 meta-IR, enabling the framework to have multiple dialects  
 126 coexisting in the same MLIR file. This approach promotes  
 127 the reuse of already defined abstractions and associated tools,  
 128 enabling intra- and inter-dialect transformations.

129 In support of the underlying algorithms and kernels used  
 130 by many machine learning frameworks (e.g., TensorFlow  
 131 and PyTorch), MLIR offers a linear algebra dialect called  
 132 *linalg* that exposes (named) operations such as convolu-  
 133 tions, matrix multiplications, and others. Operations expressed  
 134 in higher-level dialects can target *linalg* operations and  
 135 leverage all subsequent transformations supported by *linalg*

```

1 #matmul_trait = {
2   indexing_maps = [
3     affine_map<(m, n, k) -> (m, k)>, // A
4     affine_map<(m, n, k) -> (k, n)>, // B
5     affine_map<(m, n, k) -> (m, n)> // C
6   ]
7   iterator_types = [
8     "parallel", "parallel", "reduction",
9   ]
10 }
11 func.func @matmul_call(...) {
12   linalg.generic #matmul_trait
13   ins (%A, %B : memref<60x80xf32>, memref<80x72xf32>)
14   outs (%C : memref<60x72xf32>) {
15     ^bb0(%a: f32, %b: f32, %c: f32):
16       %0 = arith.mulf %a, %b : f32
17       %1 = arith.addf %c, %0 : f32
18       linalg.yield %1 : f32 }
19   return }

```

(a) Linalg Abstraction with generic operation.

```

1 func.func @matmul_call(...) {
2   // Declare constants %c0 %c1 %c4 %c60 %c72 %c80 ...
3   scf.for %m = %c0 to %c60 step %c4 { // Tiling by 4,4,4
4     scf.for %n = %c0 to %c72 step %c4 {
5       scf.for %k = %c0 to %c80 step %c4 {
6         // Grab handle for the sub-tiles:
7         %sA = memref.subview %A[%m, %k] [4, 4] [1, 1] : ...
8         %sB = memref.subview %B[%k, %n] [4, 4] [1, 1] : ...
9         %sC = memref.subview %C[%m, %n] [4, 4] [1, 1] : ...
10        // Matmul computation of a 4x4x4 tile:
11        scf.for %mm = %c0 to %c4 step %c1 {
12          scf.for %nn = %c0 to %c4 step %c1 {
13            scf.for %kk = %c0 to %c4 step %c1 {
14              %3 = memref.load %sA[%mm, %kk] : !mr4x4_0
15              %4 = memref.load %sB[%kk, %nn] : !mr4x4_1
16              %5 = memref.load %sC[%mm, %nn] : !mr4x4_1
17              %6 = arith.mulf %3, %4 : f32
18              %7 = arith.addf %5, %6 : f32
19              memref.store %7, %sC[%mm, %nn] : !mr4x4_1
20            } } } } }
21    return }

```

(b) Structured Control Flow (SCF) Abstraction with tiling.

Fig. 2: MLIR representations of a Matrix-Matrix Multiplication Operation in different abstractions.<sup>1</sup>

and lower-level abstractions. Figure 2 presents an MLIR  
 matrix-multiplication (MatMul) implementation in different  
 abstractions.<sup>1</sup> The operation is initially represented using  
 a *linalg.matmul* and subsequently undergoes conversion,  
 transformation, and lowering by the compiler. In Figure 2a-  
 L11 and L17, the *linalg.matmul* is converted into a  
*linalg.generic*. The *linalg.generic* is a core MLIR  
 operation that can represent most of the *linalg named ops*,  
 by careful selection of its *operation trait*<sup>2</sup> - *indexing\_maps*  
 (L2), *iterator\_types* (L7) -, and kernel (L24 to L27).  
 Finally, the generic operation can be converted into a tiled  
 ( $4\times 4\times 4$ ) implementation of the MatMul (Figure 2b) using  
 the structured control flow (*scf*) dialect. When supporting  
 an accelerator that can process a *MatMul<sub>4x4x4</sub>* operation<sup>3</sup>,  
 the code in Figure 2b-L11 to L19, has to be replaced by the  
 runtime library calls that drive the accelerator.

1) *MLIR Memory References*: Within MLIR, memory  
 buffers exist as N-dimensional (rank=N) memory references,  
 or *memrefs*. Our proposed AXI4MLIR DMA runtime library,  
 presented in Section III-A, supports bidirectional data move-  
 ments between *memrefs* and memory-mapped buffers (raw  
 pointers), while respecting strides, sizes, and dimensions. Ac-

<sup>1</sup>We intentionally omit some MLIR code, such as constant declarations in the form of `%cX=arith.constant X:i32`, for the sake of brevity.

<sup>2</sup>See *linalg.generic* in <https://mlir.llvm.org/docs/Dialects/Linalg>

<sup>3</sup>A 2D MatMul operation is *MatMul<sub>MxNxxK</sub>*:  $C(M,N) = A(M,K) \times B(K,N)$

158 cessing the elements of an MLIR `memref` requires accessing  
 159 the values in the equivalent C struct of Figure 3. Specializing  
 160 the code for specific sizes and strides is an important proposed  
 161 optimization to leverage spatial locality and minimize control-  
 162 flow instructions, as we will observe in Section IV.

```

1 typedef struct {
2     float *allocated; // For deallocation
3     float *aligned;   // Base address
4     size_t offset;    // Offset in # of elements
5     size_t size[N];   // One size per dim
6     size_t stride[N]; // One stride per dim
7 }

```

Fig. 3: The underlying data structure of a rank==N MLIR `memref` buffer.

### 163 B. AXI Interface

164 Efficiently using the interconnect between the CPU and  
 165 the accelerator can significantly impact the overall system  
 166 performance. As part of our framework, we consider a widely  
 167 adopted bus interface in digital electronics design deployed  
 168 on SoC and Field-Programmable Gate Array (FPGA) designs,  
 169 namely Advanced eXtensible Interface (AXI) [20]. AXI pro-  
 170 vides a flexible and scalable solution for integrating custom  
 171 accelerators into a system.

172 The AXI interface provides a simple mechanism to enable  
 173 data transfers between the CPU cores and other devices. Using  
 174 AXI, the AXI-Stream (AXI-S) interface allows the developer  
 175 to quickly transfer [25] a variable-size burst of data to and  
 176 from the accelerator in a FIFO-like manner, enabling the  
 177 accelerator to consume/store the data as needed, in a streaming  
 178 manner. Within SoCs, the CPU host code controls either a  
 179 single or multiple Direct Memory Access (DMA) engines (see  
 180 Figure 1). These engines are responsible for initiating and  
 181 handling data movement requests between the main memory  
 182 and the accelerator. Additionally, the data regions in the main  
 183 memory need to be accessible to the accelerator via the AXI-S  
 184 interface. Therefore, the host code needs to allocate input  
 185 and output memory buffers using the `mmap` function, which  
 186 guarantees that only the current process has access to the  
 187 specific regions of memory. The host code is also required  
 188 to prepare/pack the input data into the data format that the ac-  
 189 celerator requires (e.g., row-major, interleaved data elements,  
 190 etc.). Our approach within AXI4MLIR is to use *MLIR - the  
 191 compiler - to generate the host code to interface with the  
 192 accelerator*, while taking advantage of the full capabilities of  
 193 the target accelerator.

### 194 III. AXI4MLIR

195 To support efficient host code generation for AXI-based  
 196 custom accelerators, we extended the MLIR framework with  
 197 the added capabilities presented in Figure 4. After the custom  
 198 accelerator is designed and the host CPU system is selected,  
 199 the user creates a configuration file with the host CPU system  
 200 details (e.g., number and size of the caches), and with a  
 201 high-level description of the accelerator capabilities (i.e., sup-  
 202 ported operations and dimensions), the available opcodes, and  
 203 possible opcode flows ①. This information is parsed ② by  
 204 the compiler, and used to find ③ suitable `linalg.generic`

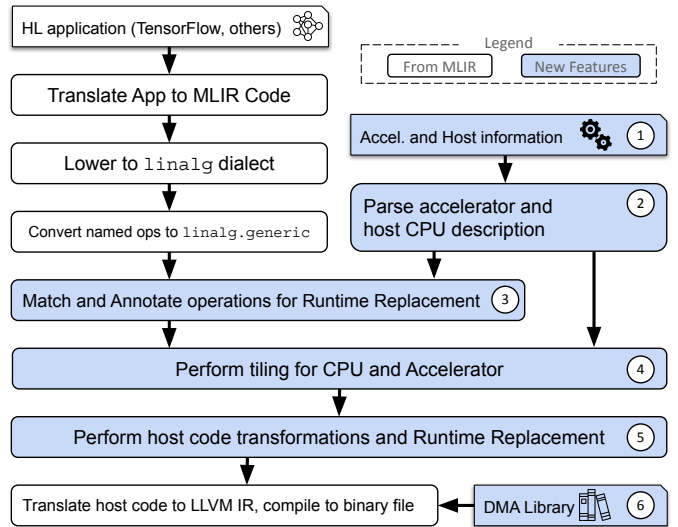


Fig. 4: AXI4MLIR Compiler Flow. The numbered elements are the contributions of this work.

operations with the desired operation traits (algorithm imple-  
 mented, previously shown in Figure 2a-L1 to L9), that can be  
 executed on the accelerator. These operations will require host-  
 accelerator driver code generation. Subsequently, with user-  
 provided information on the total size of the CPU caches, the  
 compiler transforms the code to efficiently exploit the CPU  
 memory hierarchy and the accelerator size ④, performing the  
 appropriate set of tiling transformations to leverage temporal  
 locality in the CPU caches and to map the problem on the ac-  
 celerator. In the final step, the compiler generates the runtime  
 calls ⑤ that leverage the accelerator features based on user-  
 directed dataflow description (e.g., avoiding redundant host-  
 accelerator data transfers when the algorithm and accelerator  
 functionality allows).

The following sections discuss the class of supported accel-  
 erators and the key features of our AXI4MLIR DMA library.  
 We provide details on how to describe new accelerators,  
 introducing `linalg.generic` trait extensions, a new MLIR  
 dialect that provides support for runtime call replacement of  
 opcodes and data transfers, and some key optimizations that  
 can be performed (depending on the available features of the  
 host system and the custom accelerator).

#### 227 A. The Custom AXI DMA Library

228 The AXI4MLIR DMA library ⑥ (Figure 4) exposes low-  
 229 level DMA calls working at privileged level to enable data  
 230 movement between the main memory and the accelerator. We  
 231 designed this library to be lightweight (55 bytes in size for our  
 232 target ARM SoC), so that it can be deployed on both resource-  
 233 constrained and non-constrained systems. It can also be exe-  
 234 cuted by bare-metal systems. During the compilation process,  
 235 the AXI runtime issues calls to initialize the DMA engine(s)  
 236 before entering the computation kernel of the workload. First,  
 237 a library call initializes the DMA engine, mapping memory for  
 238 the input and output buffers which act as temporary staging  
 239 buffers between the CPU and the accelerator.

240 After DMA initialization, the accelerator is accessible via  
 241 AXI-based data transfers. Any data that needs to be transferred  
 242 to the accelerator during workload execution is first copied  
 243 to a DMA input buffer. This staging copy acts as a packing  
 244 optimization (similar to [26]), contributing to an increased  
 245 cache-hit ratio during communication. Then, the AXI “send”  
 246 function call requests the DMA engine to start the data transfer  
 247 and waits for it to finish. Note that the data that is sent to the  
 248 accelerator can be either accelerator instructions or raw input  
 249 data that needs to be processed. Similarly, AXI4MLIR gener-  
 250 ates “recv” function calls to wait for computation completion  
 251 and to obtain output data from the DMA output buffer.

252 In Section III-C, Figure 9 presents the lowering of dif-  
 253 ferent high-level operations into our DMA library calls.  
 254 `copy_to_dma_region(...)` implements data movement  
 255 from a `memref` to the DMA-accessible memory region  
 256 intended for transmission to the accelerator. The `offset`  
 257 argument allows for efficient batching of different data  
 258 transfers after computing the total length and execut-  
 259 ing a single “send” operation. Appropriate offset values  
 260 prevent overwriting existing data in the DMA region.  
 261 `dma_start_send(...)` instructs the DMA engine to trans-  
 262 mit a size of  $X$  bytes to the connected accelerator, com-  
 263 mencing from a specified `offset` within the DMA space  
 264 allocated. `dma_wait_send_completion(...)` instructs the  
 265 CPU to wait for the DMA’s signal informing the transaction’s  
 266 completion. When receiving data from the accelerator, we first  
 267 have to wait for the data to be placed in the DMA-accessible  
 268 memory so it can be copied back into a `memref`.

### 269 B. Supported Accelerators

270 In matrix-multiplication and similar algorithms, the term  
 271 *stationary* refers to a slice of data that can be reused across  
 272 many iterations of an algorithm’s computation. A *stationary*  
 273 strategy attempts to maximize data reuse and minimize data  
 274 movement, which can greatly benefit accelerators that require  
 275 efficient memory accesses. We want to enable the programmer  
 276 to easily control accelerators that support *stationary* flows.

277 Next, we discuss the types of accelerators that AXI4MLIR  
 278 can support. Then we propose a standardized approach to  
 279 concisely define the class of supported accelerators in a  
 280 configuration file. Finally, we show how the AXI4MLIR parser  
 281 is able to take user-defined configurations, extract essential  
 282 attributes of the target accelerator, and populate a trait speci-  
 283 fication to guide our MLIR compiler transformations.

284 1) *Accelerator Designs*: The AXI4MLIR compiler trans-  
 285 formations support linear algebra kernels implemented as ac-  
 286 celerators using the AXI interconnect. In addition, the AXI-S  
 287 data transfers within AXI4MLIR facilitate support for accel-  
 288 erators that use a micro-ISA (Instruction Set Architecture)  
 289 with opcodes, which consist of instructions that the host-  
 290 CPU sends to the accelerator. Generally, the following three  
 291 actions are used to categorize the actions within an instruction:  
 292 *send*, *compute*, and *receive*. Any accelerator’s instructions that  
 293 require external communication (i.e., data transfers or activa-  
 294 tion/reset/configuration of accelerator compute modules) can

295 be completed by issuing a combination of these three actions.  
 296 In addition, each action can have additional meta-data (e.g.,  
 297 opcode literal, data, length, dimensions, and indexes), which  
 298 is used to guide compiler transformations during accelerator  
 299 host code generation. Further, specific traits of the accelerator  
 300 - such as internal buffer space (or accelerator tile sizes), and  
 301 data types - are supported and must be defined within the  
 302 accelerator configuration file.

```

1 { "cpu" = { "cache-levels": [32K,512K],
2           "cache-types": [data,shared] }
3   "accelerators" = [
4     { "name": ..., "version": x.x, "description": ...,
5       "dma_config" : {...}, "kernel": "linalg.matmul",
6       "accel_size": [4,4,4], "data_type": int32,
7       "dims": ["m", "n", "k"],
8       "data": { "A": [m,k], "B": [k,n], "C": [m,n] },
9       "opcode_map" : "<opcode_map string - see IV-D>",
10      "opcode_flow_map" : { "flowID01" :
11        "<opcode_flow string - see IV-D>", ... },
12      "selected_flow" : "flowID01" }]}

```

Fig. 5: Accelerator and CPU configuration file.

2) *Accelerator Configuration File*: Once an AXI-based  
 303 accelerator is fully designed, the accelerator developer can  
 304 quickly integrate it with our AXI4MLIR compiler transfor-  
 305 mations by providing *Accelerator and Host information* ①  
 306 (Figure 4) through a configuration file for the new accel-  
 307 erator and the target host system. Figure 5 shows a sample  
 308 configuration file defined in the standard JSON format. For  
 309 the accelerator, the developer must specify the accelerator’s  
 310 architectural features, e.g., supported tile sizes, data type,  
 311 and input and output data with related dimensions. Addi-  
 312 tionally, the developer should describe any micro-ISA that  
 313 the accelerator can execute. The developer should define  
 314 “opcode IDs”, captured by the “opcode\_map string”, which  
 315 are comprised of actions to describe the memory operations  
 316 and related data transfers. Finally the developer should define  
 317 the possible “opcode flow IDs” and select the desired flow  
 318 for the particular operation. The configuration file does not  
 319 capture the internal behavior of the accelerator, which has  
 320 been the focus of other works [12], [15]; instead, we seek  
 321 to optimize the communication with the accelerator. Thus the  
 322 configuration file contains information about the I/O interface  
 323 for sending data and instructions to the accelerator. Similar to  
 324 the accelerator information, the CPU information, shown in  
 325 Figure 5-L1 to L2, needs to contain basic architectural details  
 326 such as the number and size of caches.

327 3) *Configuration Parsing*: The parser implemented in ②  
 328 (Figure 4) is responsible for providing the information from  
 329 the configuration file to the MLIR IR and the AXI4MLIR  
 330 transformation passes. To this end, the *kernel* and *cache*  
 331 *information*, paired with a simple heuristic that identifies the  
 332 dimensions of the target MLIR operation, are used to schedule  
 333 tiling transformations (Figure 4 - ④) that leverage the CPU  
 334 memory hierarchy sizes and increase temporal locality of  
 335 the memory accesses. Additionally, the parser validates the  
 336 `opcode_map` and the user selected `opcode_flow`, which  
 337 are then translated into new MLIR attributes to the target  
 338

```

1 #matmul_accel_trait = {
2   dma_init_config = {           id = 0x0,
3     inputAddress = 0x42,       inputBufferSize = 0xFF00,
4     outputAddress = 0xFF42,   outputBufferSize = 0xFF00 },
5
6   // Opcodes sent once. Tokens defined in opcode_map.
7   init_opcodes = init_opcodes < (reset) >,
8
9   accel_dim = map<(m, n, k) -> (4, 4, 4)>, // Tiling
10
11  // Permutation and who can be stationary.
12  permutation_map = affine_map<(m, n, k) -> (m, k, n)>,
13
14  opcode_map = opcode_map < // Valid Opcodes
15    sA = [send_literal(0x22), send(0)],
16    sB = [send_literal(0x23), send(1)],
17    cC = [send_literal(0xF0)],
18    rC = [send_literal(0x24), recv(2)],
19    sBcCrC = [send_literal(0x25), send(1), recv(2)],
20    reset = [send_literal(0xFF)] >,
21
22  // Flow to implement. Tokens defined in opcode_map.
23  opcode_flow = opcode_flow < (sA (sBcCrC)) > // As
24  // Example of other < ((sA sB cC) rC) > // Cs
25  // valid flows < (sB sA cC rC) > // Ns
26 }

```

(a) New Attributes for Accelerator Description.

```

1 func.func @matmul_call(...) {
2   // Declare constants (loop bounds and literals): %cX, ...
3   accel.dma_init(%c0,%c66,%c65280,%c65346,%c65280) : ...
4   accel.sendLiteral(%c0xFF, %c0) : i32,i32->i32 // reset
5   // Tiling by 4,4,4
6   scf.for %m = %c0 to %c60 step %c4 { // first loop
7     scf.for %k = %c0 to %c80 step %c4 { // second loop
8       %sA = memref.subview %A[%m, %k] [4, 4] [1, 1] : ...
9       %offset0 = accel.sendLiteral(%c0x22,%c0):i32,i32->i32
10      accel.send(%sA, %offset0) : !mr4x4_0, i32 -> i32
11      scf.for %n = %c0 to %c72 step %c4 { // innermost
12        %sB = memref.subview %B[%k, %n] [4, 4] [1, 1] :
13        %sC = memref.subview %C[%m, %n] [4, 4] [1, 1] :
14        %offset1 = accel.sendLiteral(%c0x25,%c0) : ...
15        %offset2 = accel.send(%sB, %offset1) :
16          !mr4x4_0, i32 -> i32
17        accel.recv {mode="accumulate"}(%sC, %c0) :
18          !mr4x4_0, i32 -> i32
19      } } } } }
20   return

```

(b) IR to drive the MatMul accelerator with an A-stationary flow.

Fig. 6: Information added to the `linalg.generic` traits to capture accelerator behavior in MLIR and IR with `accel` operations.

339 `linalg.generic` operation trait. Their syntax and function-  
340 ality are described in Section III-C.

341 4) *Supported Systems*: Our work is focused on SoCs with  
342 accelerators connected to ARM CPUs via an AXI-S intercon-  
343 nect. AXI4MLIR seamlessly integrates with a diverse set of  
344 Xilinx platforms, though we also anticipate similar applicabil-  
345 ity to other FPGA-SoC devices. Changing the cross-compiler  
346 would allow support for other processors. Adapting our DMA  
347 library implementation to other standards would be required  
348 to support other types of interconnects. AXI4MLIR currently  
349 supports AXI-Stream accelerators, which do not communicate  
350 via direct memory requests. Thus, AXI4MLIR does not require  
351 support for host-accelerator coherence protocols, since the host  
352 manages the DMA engine transfers.

### 353 C. MLIR extensions and optimizations

354 To implement *match and annotate operations for runtime*  
355 *replacement* ③ (Figure 4), and to offload the computation  
356 onto the accelerator, we implemented passes to identify the  
357 target algorithms supported by the accelerator and extended  
358 the `linalg.generic` operation trait with additional infor-  
359 mation, as shown in Figure 6a. In particular, we introduced

```

1 opcode_dict ::=
2   "opcode_map" "<" opcode_entry ("," opcode_entry)* ">"
3 opcode_entry ::= (bare_id | string_literal) "=" opcode_list
4 opcode_list ::= "[" opcode_expr ("," opcode_expr)* "]"
5 opcode_expr ::= "send" "(" bare_id ")"
6   | "send_literal" "(" integer_literal ")"
7   | "send_dim" "(" bare_id ")"
8   | "send_idx" "(" bare_id ")"
9   | "recv" "(" bare_id ")"

```

Fig. 7: Opcode Map Syntax. A dictionary for accelerator opcodes and actions.

```

1 opcode_flow_entry ::= "opcode_flow" "<" flow_expr >
2 flow_expr ::= "(" flow_expr ")" | bare_id "(" bare_id ")"

```

Fig. 8: Opcode Flow Syntax. The sequence of opcodes to implement a specific dataflow of host-accelerator communication.

360 two new types of attributes to MLIR, `opcode_map` and  
361 `opcode_flow`, which follow the syntax described in Figure 7  
362 and Figure 8, respectively. We elaborate more on each attribute  
363 in the operation trait below. 363

### 364 Extensions to `linalg.generic` traits:

365 - `dma_init_config`: defines the parameter values used to  
366 configure a DMA engine associated with a specific accelerator.  
367 If multiple or different accelerators are present, they would  
368 have different values in this field. Figure 6a-L2 to L4 show  
369 the available parameters. The code generated for the DMA  
370 initialization is executed by the CPU only once per application. 370

371 - `init_opcodes`: defines a flow of opcodes that should be  
372 sent to initialize or reset the accelerator for a new kernel  
373 execution. During application runtime, these opcodes are sent  
374  $N$  times, where  $N$  is the number of kernels in an application  
375 that can be mapped onto the custom accelerator. In Figure 6a-  
376 L7, we define that the reset opcode must be included to support  
377 the described accelerator. The opcode's functionality is derived  
378 from the `opcode_map` parameter below. 378

379 - `accel_dim`: defines the size of the accelerator for each  
380 dimension of the implemented algorithm. Figure 6a-L9 shows  
381 an example, specifying that the *accelerator* supports a tiled  
382 `MatMul4x4x4` version of the implemented algorithm. 382

383 - `permutation_map`: defines the order in which nested loops  
384 execute. In Figure 6a-L12, we switch the order of the two  
385 innermost loops, potentially enabling the data structure that  
386 uses  $[m,k]$  indices to be stationary, as the other data structures  
387 are streamed in/out of the accelerator. In our MatMul example  
388 (Figure 2b), this enables an A stationary dataflow (Figure 6b). 388

389 - `opcode_map`: describes accelerator opcodes as key-value  
390 pairs. Following the syntax scheme shown in Figure 7, the  
391 key, or *opcode\_entry*, is an identifier that maps to a list of  
392 actions, or *opcode\_list*, which represents sequential memory  
393 operations that have to be performed to drive the accel-  
394 erator. Each action, or *opcode\_expr* (`send`, `send_literal`,  
395 `send_dim`, `send_idx`, `recv`), implements different types  
396 of copies to/from the DMA memory-mapped region. The  
397 `send` and `recv` actions take an input. The input is a num-  
398 ber that is used to represent one of the arguments to the  
399 `linalg.generic` operation, e.g., 0, 1, or 2 would map to  
400 A, B, or C, respectively, in the MatMul example (Figure 2a-  
401 L12-13). During code generation, this information is used  
402 to copy the needed tile to the memory-mapped region. For



Fig. 9: Semantics and lowering of accel operations.

example, Figure 6a-L15 shows an opcode with identifier “sA” that issues copies *to* the accelerator for the *literal* value 0x22 and then for the *data* associated with the tile of argument 0. Furthermore, `send_dim` and `send_idx` can be used to send tile dimensions or tile indices, which could be used to drive more complex accelerators. Subsequent text will refer to an *opcode\_entry*, such as “sA”, simply as *opcode*.

– *opcode\_flow*: represents valid opcode/data transfer *flows* and respects the syntax scheme shown in Figure 8. Figure 6a-L23 shows an example, which defines an *input A stationary* (associated with argument 0) valid flow implemented with two opcodes, using the identifiers defined in the `opcode_map`. Additional valid examples for *output C stationary* and *nothing stationary* flows are shown in lines 24 and 25 of Figure 6a. The information in `opcode_flow` is parsed and the set of parentheses is understood as a proxy to specify multiple scopes for sequential or nested *for* loops in the algorithm. Following this flow, logic related to “sA” would be transmitted inside of the second loop (Figure 6b-L8 to L10), and logic related to “sBcCrC” would appear in the innermost loop (Figure 6b-L12 to L18). Suppose the user decides to forego the opportunity to specify *input A* as stationary, then the opcode flow could become “(sA sB cC rC)”, and all communication driver logic would be generated in the innermost loop.

**The accel dialect:** Before generating function calls for *runtime replacement* to the DMA runtime library (described in Section III-A), we perform *host code transformations* (5) (Figure 4) by lowering the `linalg.generic` operation, with the proposed `trait`, to standard MLIR dialects (`scf`, `arith`, `memref`) and a new dialect that we call `accel`. Operations in the `accel` dialect abstract host-accelerator transactions, such as initialization, memory transfers, and synchronization. Figure 9 presents the core `accel` operations and their semantics, providing examples of how these operations map onto our custom AXI DMA library calls. Additionally, Figure 6b shows how the `accel` operations are used in our MatMul example.

Note that it is easier to perform analysis and transformations

of operations when they are expressed in our `accel` dialect, as opposed to using a lower-level abstraction. With lower-level abstractions such as `llvm`, function calls and additional logic have already been exposed: additional instructions must be present in the IR to implement buffer slicing, size/offset calculations, and function calls to copy data to/from the DMA regions. Performing analysis and transformations in the `llvm` abstraction is more challenging, as traversal of control flow blocks and LLVM instructions are necessary. Instead, operations in the intermediate `accel` dialect encode the relevant information, and are easily relocated during transformation passes, respecting dependencies without requiring complex compiler analysis. This approach facilitates implementing communication flows that consider one of the data structures to be stationary by simply hoisting the `accel` operations up to the right loop nest level, while considering the flow patterns. Finally, the `accel` dialect provides an intermediate step before runtime call replacement. In this work we target our AXI DMA runtime library described in Section III-A, but further extensions could implement the transformation of `accel` operations into other runtime libraries such as OpenCL [27] or SYCL [28], which are commonly used to interface with SoC FPGA accelerators.

#### IV. EXPERIMENTS AND RESULTS

To evaluate AXI4MLIR, we use a PYNQ-Z2 board that includes a Zynq-7000 SoC with a dual-core ARM Cortex-A9 CPU (650 MHz), and a library of tile-based accelerators derived from SECDA-TFLite [29] implemented with AXI-S interface and opcodes with a micro-ISA. For workloads, we target a suite of kernels covering a range of dimensions, as well as an end-to-end machine learning application. We leverage hand-written baselines, which we discuss in Section IV-A. Section IV-B evaluates accelerators implementing MatMul, comparing inference performance against a hand-written baseline, identifying potential bottlenecks, and showcasing the benefits of our optimized dataflows. Section IV-C highlights the value of AXI4MLIR by demonstrating how to handle accelerators with configurable parameters such as tile sizes and dataflows. We showcase how to use AXI4MLIR with a convolution-based accelerator in Section IV-D. Finally, Section IV-E shows how AXI4MLIR can work in the context of a complete application, evaluating the TinyBERT model [30].

##### A. Hand-written Baselines

The next experiments employ hand-written optimized driver code derived from the SECDA-TFLite accelerator toolkit [29] to establish performance baselines. SECDA-TFLite presents a state-of-the-art toolchain and methodology for HW/SW co-design of embedded machine learning accelerators targeting FPGA SoC devices. With host-driver code written in C++, these manual baselines will be labeled as `cpp_MANUAL`. All baselines are implemented with various tiling strategies, with no additional data transfer overheads and with the fewest number of data transfer calls for the selected dataflow.

493 **B. Matrix-Multiplication Experiments**

494 The tile-based accelerators used here resemble vector MAC  
 495 engines [31], [32], [33], [34] implementing MatMul algo-  
 496 rithms. They vary in input/output buffer size and supported  
 497 dataflow. From the CPU-host perspective, some of them can  
 498 support varying degrees of data reuse when the appropri-  
 499 ate opcode stream drives the accelerator. Table I presents  
 500 a short summary of their functionality, where *size* stands  
 501 for the supported tile size of the accelerator. For example,  
 502  $v1_4$  is a  $\text{MatMul}_{4 \times 4 \times 4}$  accelerator that does not support  
 503 data reuse and only supports  $tM, tN, tK == 4, 4, 4$  tiles.  
 504 For  $v1_4$ , AXI4MLIR will tile the algorithm’s loops in the  
 505 host code, taking into account the accelerator size of 4 and  
 506 all the data movement will happen in the innermost loop  
 507 - “opcode\_flow <(sA sB cCrC)>”. For  $v2_8$ , AXI4MLIR  
 508 will tile the computation by 8 and generate code to maximize  
 509 the reuse of one of the inputs. In  $v2$ , a stationary (As) is  
 510 implemented with opcode\_flow <(sA (sB cCrC))>.

511 Accelerators  $v3$  and  $v4$  can also reuse their output data  
 512 structures. Accelerator  $v4$ , marked with *flex size*, supports  
 513 computations of non-square tiles, i.e.,  $v4_{16}$  can process a  
 514  $\text{MatMul}$  of  $tM, tN, tK = 32, 16, 64$ , as long as  $tM, tN, tK$   
 515 are divisible by 16 and fit in the accelerator’s memory. All  
 516 accelerators were implemented using HLS pipelining and  
 517 unrolling to maximize the number of internal processing  
 518 elements instantiated and their arithmetic throughput. The  
 519 last column of Table I reports throughput (OPS/cycle) for  
 520 each accelerator, highlighting that many arithmetic operations  
 521 are executed in parallel at each cycle. Different types of  
 522 accelerators with the same size have the same throughput, and  
 523 accelerators with bigger sizes provide higher throughput. All  
 524 bar graphs presented in this section represent the average of  
 525 5 independent runs with the same configuration.

526 **Accelerator relevance.** In order to evaluate the performance  
 527 of the accelerators defined in Table I, we conducted ex-  
 528 periments to compare the runtime of the CPU execution  
 529 (*mlir\_CPU*) against the manual C++ implementation (referred  
 530 to as *cpp* for short) of the driver code using the accelera-  
 531 tors. The task clock was used as a metric to measure the  
 532 execution time of the benchmarks. We present the results of  
 533 the experiments in Figure 10, which plots the task clock on  
 534 the y-axis (smaller is better) and only includes the “Nothing  
 535 Stationary flow”, which means that the data transfers happen  
 536 in the innermost loop.

537 Looking at Figure 10, we can see that the accelerator offload  
 538 only becomes relevant (i.e., executes faster than the CPU) for  
 539 problems with  $\text{dims} \geq 64$ , where  $\text{dims} = M = N = K$ .  
 540 For problems with smaller dimensions, CPU execution will be  
 541 faster than the accelerator. In addition, the results in Figure 10  
 542 suggest that accelerators only become relevant if  $\text{accel\_size} =$   
 543  $tM = tN = tK \geq 8$ . For smaller accelerator sizes, the CPU  
 544 execution is faster than the accelerator.

545 These observations suggest that the performance benefits of  
 546 using the accelerators are limited for ranges of problem sizes  
 547 and accelerator sizes. Therefore, it is important to carefully

TABLE I: Accelerators used in the experiments. Synthesized with AMD/Xil-  
 inx Vitis at 200MHz.

Type	Possible Reuse	Opcode(s)	Configurations (Size, OPs/Cycle)
$v1_{\text{size}}$	Nothing	sAsBcCrC	(4, 10)
$v2_{\text{size}}$	Inputs	sA, sB, cCrC	(8, 60)
$v3_{\text{size}}$	Ins/Out	sA, sB, cC, rC	(16, 112)
$v4_{\text{size}}$	Ins/Out (flex size)	sA, sB, cC, rC	

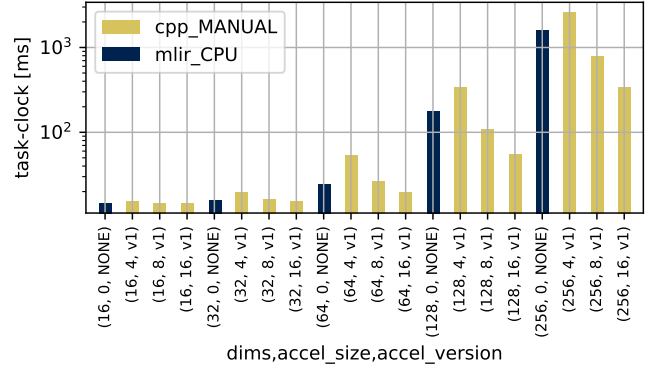


Fig. 10: Runtime characterization CPU vs. Accelerator execution for Matrix  
 Multiplication problems. Note how an accelerator only becomes relevant for  
 problems with  $\text{dims} \geq 64$  and  $\text{accel\_size} \geq 8$ .

548 choose the appropriate accelerator configuration for a given  
 549 problem to achieve the best performance. Consequently, for  
 550 the next experiments we will limit our focus to problems with  
 551  $\text{dims} \geq 64$  and accelerators with  $\text{accel\_size} \geq 8$ .

552 **AXI4MLIR generated vs. Manual implementation.**

553 AXI4MLIR provides several benefits. First, our passes au-  
 554 tomatically tile data mapped to the CPU memory hierarchy,  
 555 leveraging spatial and temporal locality. The second benefit is  
 556 the ability to automatically generate specific flows, such as the  
 557 Nothing Stationary (Ns) flow, which can be tedious and error-  
 558 prone when done manually. Additionally, AXI4MLIR provides  
 559 an efficient path to flow strategies that can potentially improve  
 560 performance, such as input A or B stationary (As, Bs) flows.  
 561 Figure 11 presents these results.

562 First, we compare the differences in execution time between  
 563 a *manual implementation* (see Section IV-A) of an Ns flow  
 564 strategy and an *AXI4MLIR generated* Ns flow strategy, repre-  
 565 sented by the first two bars in each group of bars in Figure 11.  
 566 The remaining bars in each group of bars show results for  
 567 automatically generated flow strategies, with As and Bs for  $v2$   
 568 accelerators and As, Bs, and Cs for  $v3$  accelerators. Looking  
 569 at Figure 11 we see that some flows, especially Cs, provide  
 570 improvements. To achieve this, the user simply has to encode  
 571 the information for Cs (or other flows) during compilation. For  
 572 example, we can encode Cs using the opcode\_flow previously  
 573 presented in Figure 6a-L25 in the the operation’s trait.

574 Next, in Figure 11, we focus on the results with the  
 575  $v3$  accelerator. Here, we see that AXI4MLIR generated Cs  
 576 performs better than the manually generated Ns, although the  
 577 other flows are not performing as expected. First, we would  
 578 expect the performance of AXI4MLIR generated Ns to have  
 579 similar/closer task clock performance than manual Ns. And  
 580 second, we would also expect As and Bs flows to always  
 581 outperform Ns due to the degree of reuse, as they copy less



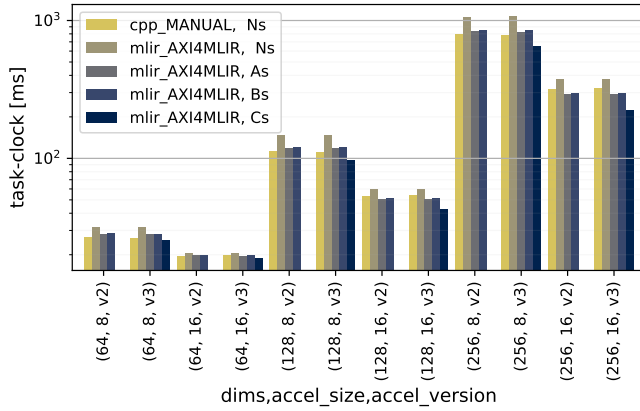


Fig. 11: Runtime results on Matrix Multiplication kernels. Manual implementation of Ns flow vs. AXI4MLIR generated driver code for different flow strategies, Ns, As, Bs, Cs. All bar groups follow similar trends. Ns, As, and Bs **bottlenecks are analysed and addressed** in following experiments.

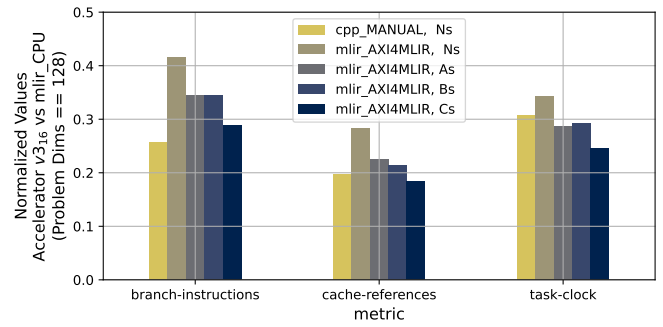
582 data and can keep the accelerator better utilized. Hence, this  
 583 first implementation has room for improvement and, in the  
 584 following experiment, we *identified and fixed the bottlenecks*  
 585 by analyzing performance counters and implementing opti-  
 586 mizations that specialize memory copies.

### 587 Identifying bottlenecks & improving AXI4MLIR codegen.

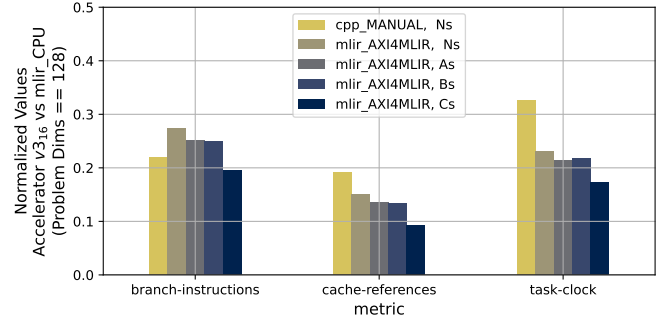
588 Next, we identify performance bottlenecks in AXI4MLIR  
 589 generated copies and improve upon them to enhance the  
 590 performance of the workloads when using the custom hard-  
 591 ware accelerators. Specifically, the experiment compares the  
 592 performance of manually implemented host-accelerator driver  
 593 code with AXI4MLIR generated code for Ns, As, Bs, and Cs  
 594 flows in terms of branch-instructions, cache reference counters,  
 595 and the task clock. These metrics were obtained using the  
 596 perf tool [35] to profile the application and retrieve counters  
 597 for CPU perf\_events over 5 runs.

598 Figure 12a shows branch instructions, cache reference coun-  
 599 ters, and the task clock for  $dims == 128$ , for the  $v_{316}$   
 600 accelerator that supports input and output stationary flows. The  
 601 trends are similar to other problem and accelerator sizes. Our  
 602 results are normalized to the same counters collected on a  
 603 CPU-only execution of the same problem size. In each group  
 604 we show results for AXI4MLIR automatically generated code  
 605 for Ns, As, Bs, Cs flows, and compare against manual imple-  
 606 mentations (first bar of a group) for copying the necessary data  
 607 through the DMA memory-mapped region. MLIR applications  
 608 have to consider MLIR memory references (presented in  
 609 Section II-A1), but manual implementations use bare C-arrays.  
 610 To support generality, MLIR copies between MemRef and  
 611 the raw array (DMA buffer region) are implemented with a  
 612 recursive call, loading and storing one element at a time. This  
 613 is necessary to support  $rank = N$  MemRefs, where strides in  
 614 all dimensions are different from 1.

615 In order to address this issue, we implemented an opti-  
 616 mization for when  $strides[N - 1] == 1$  (i.e., elements in  
 617  $N - 1$  dimension are adjacent to each other in memory) and  
 618 specialized MemRef copies for some known rank sizes, such  
 619 as  $rank == 2$ . For this scenario, we leverage the spatial



(a) Without the MemRef-DMA buffer copy optimization. Generated host-accelerator code has overheads if not specialized.



(b) With MemRef-DMA buffer copy optimization. AXI4MLIR improves accelerator performance over manual Ns implementation.

Fig. 12: Cache, branch, and runtime metrics of different tools and strategies using  $v_{316}$  accelerator with problem size ( $dims == 128$ ). Normalized values against CPU (without accelerator) executions of same problem size.

620 locality and implement the copy not with individual load and  
 621 store instructions, but by calling `std::memcpy(src, dst, size)`.  
 622 When compiling this function for our platform, the  
 623 compiler will inline the assembly, implementing a vectorized  
 624 copy that improves the performance of the copy operation.  
 625 The implications of this optimization are twofold. First, it  
 626 reduces the number of branch references because there is no  
 627 need for branching to handle non-unitary strides or to iterate  
 628 over an arbitrary number of dimensions, resulting in better  
 629 control flow and branch prediction. Second, the vectorized  
 630 code reduces the number of cache references because the data  
 631 is accessed sequentially in memory. Therefore, there will only  
 632 be two cache reference to fetch the cache line containing the  
 633 requested data, and subsequent loads within the same cache  
 634 line will not require additional cache references as they are  
 635 read from the vector VFP registers [36]. The results for this  
 636 optimization are presented in Figure 12b.

637 After incorporating this optimization, the AXI4MLIR gen-  
 638 erated driver code executed faster on all accelerators as  
 639 compared to their respective manual implementations. In Fig-  
 640 ure 13, we compare AXI4MLIR against manual implemen-  
 641 tations for Ns, As, Bs, and Cs, and found that the compiled  
 642 generated driver code provided by AXI4MLIR is consistently  
 643 faster ( $1.18\times$  average speedup and  $1.65\times$  max speedup),  
 644 thanks to its ability to leverage proper tiling for the CPU's  
 645 memory hierarchy, resulting in a 10% average and 56% max  
 646 reduction in cache references.

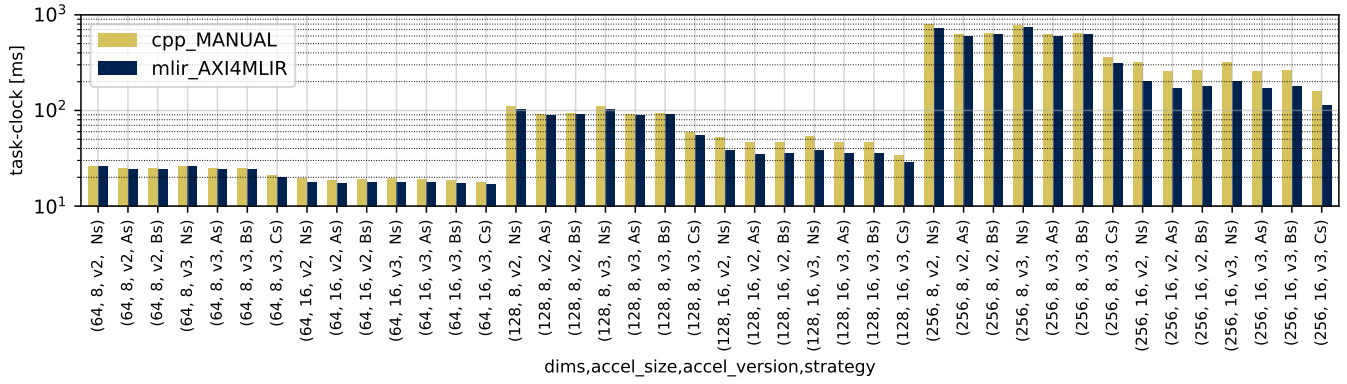


Fig. 13: Runtime comparison of manual implementation of driver code and AXI4MLIR generated. Each set of two bars have a matching Accelerator Type, Accelerator Size, and Flow Strategy (Ns, As, Bs, Cs). AXI4MLIR is better in all cases.

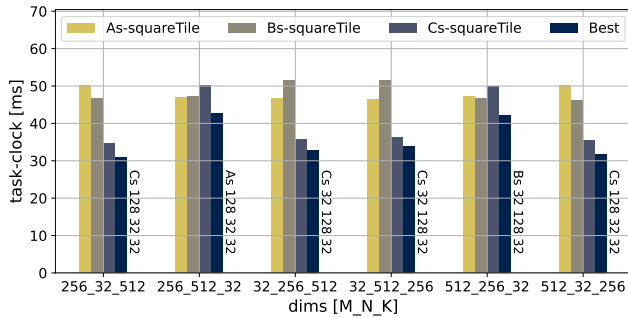


Fig. 14: MatMul problem permutations (v4 accelerator) for different strategies. For the “Best” strategies we annotate the chosen flow and tiling values.

### C. Matrix-Multiplication with flexible sizes

Runtime configurable accelerators allow for fine-grained hardware tuning for specific problems. With AXI4MLIR, we can generate host code to configure and optimize flexible accelerators for the target problem. To demonstrate this capability, we evaluate multiple permutations of a MatMul problem on the v4 accelerator. The v4 accelerator supports multiple dataflow strategies and adjustable tile sizes for its  $tM$ ,  $tN$ , and  $tK$  dimensions. The intuition is that scientific and machine learning workloads present problem sizes with different values for each dimension, sometimes resulting in tall/skinny matrices during execution. Tiling the problem in the accelerator with different dimensions for  $tM$ ,  $tN$ , and  $tK$ , and selecting the appropriate flow strategy can be beneficial for the application.

When using AXI4MLIR, a developer is *not* limited to one configuration of an accelerator. Based on the user’s knowledge of the application, AXI4MLIR can automatically generate the driver for accelerators with adjustable dimensions. This flexibility allows for a more thorough exploration of the design space, enabling the developer to find the best sizes for  $tM$ ,  $tN$ ,  $tK$ , and the best flow strategy for each problem instance.

In Figure 14, we compare four different heuristics and use them to choose the best tiling and dataflow configuration for a MatMul problem. We evaluate performance in terms of execution time. We profile the problem with  $M, N$ , and  $K$  dimensions permuted from the following values: [32, 256, 512]. Hence, the theoretical minimum number of multiply-accumulate operations required for all permutations

is the same. Here, the *As-squareTile*, *Bs-squareTile*, and *Cs-squareTile* heuristics try to find the best configuration to reduce the total memory access count given the constraint of tiling the MatMul with square tiles (i.e.,  $tM = tN = tK = T$ ), with A, B, and C stationary dataflow, respectively. The fourth heuristic, *Best*, chooses between all dataflows and flexible tiling options, only sharing the choice of the accelerator. In Figure 14 we annotate the “Best” configuration found for each problem.

**Square tiling.** We observe that as we change the problem permutation, the best flow between *As-squareTile*, *Bs-squareTile*, and *Cs-squareTile* tiling strategies changes. The best flow depends on the problem shape, the size, and the available accelerator buffer space.  $T = 32$  was selected for all square flows because it is the biggest value, so the tiles fit inside the accelerator’s internal memory.

**Flexible tiling.** The *Best* heuristic, selected from non-square strategies, outperforms square tiling by leveraging flexible tiling sizes. AXI4MLIR can generate code to utilize larger tile sizes in various dimensions, taking advantage of the v4 accelerator’s unrestricted tiling factors and improving the accelerator’s internal memory utilization.

**Configurations.** Manually implementing all configurations’ driver code for even a simple accelerator such as v4 is very time-consuming. AXI4MLIR can quickly generate hostcode for configurable accelerators easily, enabling the developer to specify an accelerator configuration per problem instance.

### D. Convolution

We show the flexibility of AXI4MLIR by generating driver code for a convolution-based accelerator executing different problems sizes. This accelerator supports varying input channel ( $iC$ ) and filter ( $fHW$ ) sizes, computing one output slice (all elements in one output channel -  $oC$ ) per iteration. To orchestrate the execution, multiple instructions have to be sent to the accelerator. This orchestration is achieved by compiling the driver code derived from the MLIR `accel` code (Figure 15b). The `accel` code is generated after a transformation pass takes into account the attributes shown in Figure 15a and MLIR’s `linalg.conv_2d_nchw_fchw` operations. Note that if the convolution operation has  $iC$ ,  $fH$ ,  $fW$  dimensions that are smaller than the dimensions in `accel_dim`, no tiling

```

1 accel_dim = map<(B,H,W, iC,oC,fH,fW) ->
2   (0,0,0,256, 1, 3, 3)>, // Tiling
3 opcode_map<
4   sIc0=[send_literal(70), send(0)], // send 3D input window
5   sF=[send_literal(1), send(1)], // and compute
6   rO=[send_literal(8), recv(2)], // recv 2D output slice
7   rst=[send_literal(32), send_dim(1,3), // set filter size
8     send_literal(16), send_dim(0,1)]> // set iC size
9 opcode_flow <(sF (sIc0) rO)> // filter+output stationary
10 init_opcodes <(rst)>

```

(a) Opcode Map and Flow for Conv2D accelerator.

```

1 func.func @conv_call(...) {
2   // With %I: !mrI_1_256_7; %W: !mrW_64_256_3_3
3   // and %O: !mrO_1_64_5_5
4   // Declare constants (loop bounds and literals): %cX, ...
5   accel.dma_init(%c0,%c66,%c65280,%c65346,%c65280) : ...
6   accel.sendLiteral(%c32, %c0) : i32,i32->i32 // send inst
7   accel.sendDim(%W,%c3,%c0) : !mrW,i32,i32->i32 // send %fH
8   accel.sendLiteral(%c16, %c0) : i32,i32->i32 // send inst
9   accel.sendDim(%I,%c1,%c0) : !mrI,i32,i32->i32 // send %iC
10
11  // Tile dims by (B,H,W,iC,oC,fH,fW) -> (-,-,-,256,1,3,3)
12  scf.for %b = %c0 to %c1 step %c1 { // B loop
13    scf.for %oc = %c0 to %c64 step %c1 { // OC loop
14      %sW = memref.subview %W[%oc,0,0,0][1,256,3,3] ...
15      %offset0 = accel.sendLiteral(%c1, %c0) : i32,i32->i32
16      %offset1 = accel.send(%sW, %offset0) :
17        !mrSubWx256x3x3, i32 -> i32
18      scf.for %oh = %c0 to %c5 step %c1 { // OH loop
19        scf.for %ow = %c0 to %c5 step %c1 { // OW loop
20          %xoffset = ... // index calculation
21          %yoffset = ... // index calculation
22          %sI = memref.subview %I[0,0,%xoffset,%yoffset]
23            [1,256,3,3] ...
24          %offset2 = accel.sendLiteral(%c70, %c0) : ...
25          %offset3 = accel.send(%sI, %offset2) :
26            !mrSubIx256x3x3, i32 -> i32
27          // inner product of sW and sI computed in HW
28        }
29        %sO = memref.subview %O[0,%oc,0,0] [1,1,5,5] ...
30        %offset4 = accel.sendLiteral(%c8, %c0) : ...
31        accel.recv {mode="accumulate"}(%sO, %c0) :
32          !mrSubO_5x5_0, i32 -> i32
33      } } } return }

```

(b) IR to drive the Conv2D accelerator with an output-stationary flow.

Fig. 15: Information added to the linalg.generic traits to capture convolution accelerator behavior in MLIR and IR with accel operations.

will be performed across these dimensions. In the convolution example (Figure 15), upon accelerator reset, we use `send_dim(1, 3)` to send to the accelerator the filter size as the dimension ‘3’ of data structure ‘1’ (i.e., the filter), and we use `send_dim(0, 1)` to send the input channel size as the dimension ‘1’ of the data structure ‘0’ (i.e., the input).

We evaluate the performance of AXI4MLIR during the execution of all convolution layers of ResNet18 [37]. Figure 16 presents performance metrics normalized to the runtime of layer-specific manual C++ driver code. The results observed here present similar trends as observed in the MatMul experiments. Only one layer (`56_64_1_128_2`) presented a 10% slowdown, contrary to previous trends. The slowdown happened because `fHW` (1) and `iC` (64) were too small, and the overhead of dealing with small MemRefs was not overcome since we could not leverage the *strided copy optimization* presented in Section IV-B. Smaller AXI4MLIR speedups are observed every time that `fHW`==1. That said, AXI4MLIR achieves better runtime performance on 10 out of 11 ResNet18 layers, with 1.28× and 1.54× average and max speedup, respectively, thanks to the improved CPU cache performance.

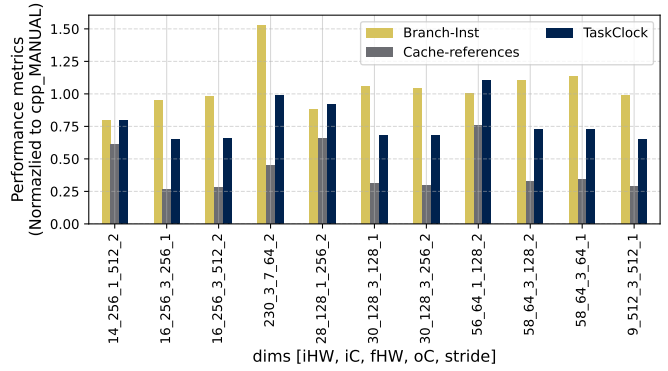


Fig. 16: ResNet18 convolution layers: AXI4MLIR vs. Manual.

## E. End-To-End Analysis

Finally, we evaluate AXI4MLIR when compiling a natural language processing model to co-execute on both the CPU and the v4<sub>16</sub> accelerator. We benchmark TinyBERT [30], a compact transformer [38] model for Masked Language Modeling and Next Sentence Prediction targeted at mobile and embedded devices. We translate TinyBERT to MLIR IR using TorchMLIR [39] and compare the inference performance of CPU execution (using -O3 during compilation) against co-execution using the “Ns” offloading approach and the “Best” approach, which employs the heuristics presented in Section IV-C.

As we can see in Figure 17, AXI4MLIR achieves a 3.4× speedup in end-to-end execution, with an 18.4× speedup in the accelerated MatMul layers that represent 75% of the original CPU runtime. This experiment showcases how AXI4MLIR can be used during evaluation and optimization of natural language processing models on embedded devices. Our study highlights that developers can easily co-design the accelerators when targeting full workloads, which enables efficient exploration and utilization of both CPU and accelerator resources.

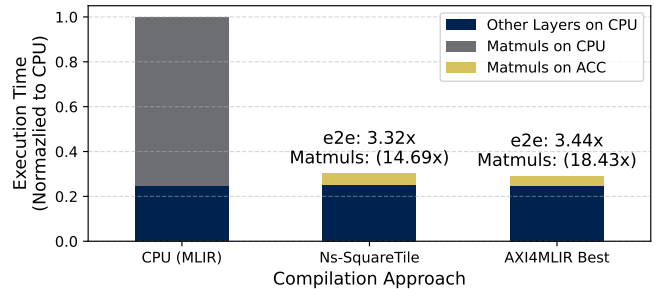


Fig. 17: Execution time of the TinyBERT model with `batch_size` == 2. Each bar represents a compilation strategy. Speedups for end-to-end (e2e) and for accelerated MatMul layers are shown as annotations.

## V. RELATED WORK

Prior studies [40], [17], [11], [16], [13], [41], [42], [43] have proposed new accelerator designs or presented new methodologies to generate flexible accelerators for a wide range of algorithms. However, these approaches fall short in providing insights into how host-to-accelerator transfers should be performed. Most of these tools assume that the required data

763 is already placed in the accelerator’s internal buffers. There  
764 have also been efforts to support hardware/software co-design  
765 of an accelerator for an application [18], [19], [29]. However,  
766 these implementations adopt a simple offload model, where  
767 execution of the kernel code is simply *replaced* by runtime  
768 calls that copy the data to-and-from the accelerator, *without*  
769 *considering the host memory hierarchy or accelerator features*,  
770 which would require manual driver code modifications.

771 HeteroFlow [44], an FPGA accelerator programming model,  
772 decouples algorithm specification from data placement op-  
773 timization using a new primitive “.to()”. This approach  
774 exposes data placement specification at various granularities,  
775 achieving efficient code generation while matching optimized  
776 manual HLS designs. HeteroFlow does not support arbitrary  
777 custom accelerators, as it is limited to accelerators co-designed  
778 with its framework (extended HeteroCL [45]). It also requires  
779 the new primitive to be used while describing the algorithm  
780 in Python, imposing manual application modification. Unlike  
781 HeteroFlow, AXI4MLIR utilizes MLIR to target languages  
782 employing `linalg.generic` operations during compilation,  
783 eliminating the need for manual transformation.

784 Several other studies have addressed the challenge of effi-  
785 ciently mapping algorithms and their loops onto accelera-  
786 tors through operation scheduling. Notably, Interstellar [46],  
787 DMazeRunner [47], and PolySA [48] delve into more versatile  
788 loop structures by adopting diverse loop representations for  
789 DNN layers. CoSA [49] and Vaidya et al. [50] tackle the  
790 generation of execution schedules for DNNs in a time-efficient  
791 manner, leveraging constrained optimization solvers. Self-  
792 tuning algorithms have also been employed in addressing the  
793 scheduling problem. Approaches like ConfuciuX [51], Flex-  
794 Tensor [52], AutoTVM [53], and AnsoR [54] utilize machine  
795 learning algorithms. Furthermore, Flexer [55] employs an out-  
796 of-order scheduling technique, unbound by loop order, which  
797 orchestrates operations based on a comprehensive analysis of  
798 the data-flow graph of a given layer. Some of these works  
799 are tailored to a specific type of accelerator or algorithm.  
800 In addition, these works primarily focus on scheduling as-  
801 pects, which AXI4MLIR currently lacks as a component.  
802 Nonetheless, these scheduling techniques could potentially  
803 complement AXI4MLIR’s attributes and transformations to  
804 enhance the overall accelerator integration process.

805 The Pattern Description Language (PDL) and Transform  
806 MLIR dialects [56] offer productive ways for expressing IR  
807 transformations and could be leveraged to implement simi-  
808 lar functionality as provided by AXI4MLIR. However, PDL  
809 cannot currently identify patterns in nested MLIR regions.  
810 Additionally, the transform dialect focuses on scheduling  
811 linear algebra transformations but requires extensions for  
812 runtime call generation targeting accelerators and dataflows.  
813 In contrast, AXI4MLIR’s `opcode_map` and `opcode_flow`  
814 extensions enable flexible automatic driver code generation  
815 for custom accelerators. Future work may involve integrating  
816 AXI4MLIR passes as Transform operations and using PDL  
817 to identify operation sequences for transformation, potentially  
818 supporting fusing operations for custom accelerator execution.

Host code generation transforms accel operations into DMA  
library calls. To facilitate further optimizations leveraging the  
MLIR infrastructure, users can modify these transformation  
passes while applying optimizations such as double buffering,  
building on our infrastructure supporting non-blocking trans-  
fers and transfer completion checks. Our ongoing work will  
introduce a new attribute to select inputs/outputs for double  
buffering. This aligns with MLIR’s capability to modify and  
add passes to the transformation pipeline. For further effi-  
ciency, coalescing transfer requests are essential; future work  
will implement a transformation that consolidates multiple  
`start_send` calls into a single call after data preparation,  
thus reducing the need for multiple `wait_send` calls, which  
incur higher CPU accelerator-DMA synchronization costs.

## VI. CONCLUSION

In this paper we presented AXI4MLIR, an extension to the  
MLIR compiler framework to describe AXI-based accelerators  
with a range of features including accelerator opcodes. We  
implemented attribute extensions and compiler transformations  
to describe and automatically generate host code that can  
leverage different flows of flexible accelerators, allowing us  
to break away from simple offload HW/SW co-design models.  
After implementing data staging and accessing optimizations  
during communication, our results show that AXI4MLIR is  
effective in generating host code that efficiently uses CPU  
resources and accelerator features. This allows for measurable  
runtime improvements versus manual implementations for all  
tested accelerators, while providing automation and conveni-  
ence during the co-design cycle. Finally, our user-driven host  
code generation is entirely automated, providing a significant  
advantage in terms of productivity and maintainability, spe-  
cially during the early stages of the co-design process.

## REFERENCES

- [1] J. Hennessy and D. Patterson, “A new golden age for computer archite-  
cture: Domain-specific hardware/software co-design, enhanced security,  
open instruction sets, and agile chip development,” in *2018 ACM/IEEE  
45th Annual International Symposium on Computer Architecture (ISCA)*.  
Los Angeles, CA, USA: IEEE, 2018, pp. 27–29.
- [2] H. Shabani, A. Singh, B. Youhana, and X. Guo, “Hirac: A hierarchical  
accelerator with sorting-based packing for spgemms in dnn applica-  
tions,” in *2023 IEEE International Symposium on High-Performance  
Computer Architecture (HPCA)*. Montreal, QC, Canada: IEEE, 2023,  
pp. 247–258.
- [3] B. Kim, S. Li, and H. Li, “Inca: Input-stationary dataflow at outside-  
the-box thinking about deep learning accelerators,” in *2023 IEEE  
International Symposium on High-Performance Computer Architecture  
(HPCA)*. Montreal, QC, Canada: IEEE, 2023, pp. 29–41.
- [4] J. Zhao, Y. Yang, Y. Zhang, X. Liao, L. Gu, L. He, B. He,  
H. Jin, H. Liu, X. Jiang, and H. Yu, “Tdgraph: A topology-driven  
accelerator for high-performance streaming graph processing,” in  
*Proceedings of the 49th Annual International Symposium on Computer  
Architecture*, ser. ISCA ’22. New York, NY, USA: Association  
for Computing Machinery, 2022, p. 116–129. [Online]. Available:  
<https://doi.org/10.1145/3470496.3527409>
- [5] S. Hsia, U. Gupta, B. Acun, N. Ardalani, P. Zhong, G.-Y. Wei,  
D. Brooks, and C.-J. Wu, “Mp-rec: Hardware-software co-design to  
enable multi-path recommendation,” in *Proceedings of the 28th ACM  
International Conference on Architectural Support for Programming  
Languages and Operating Systems, Volume 3*, ser. ASPLOS 2023.  
New York, NY, USA: Association for Computing Machinery, 2023, p.  
449–465. [Online]. Available: <https://doi.org/10.1145/3582016.3582068>

- [6] S. Zheng, R. Chen, A. Wei, Y. Jin, Q. Han, L. Lu, B. Wu, X. Li, S. Yan, and Y. Liang, "Amos: Enabling automatic mapping for tensor computations on spatial accelerators with hardware abstraction," in *Proceedings of the 49th Annual International Symposium on Computer Architecture*, ser. ISCA '22. New York, NY, USA: Association for Computing Machinery, 2022, p. 874–887. [Online]. Available: <https://doi.org/10.1145/3470496.3527440>
- [7] F. Muñoz Martínez, R. Garg, M. Pellauer, J. L. Abellán, M. E. Acacio, and T. Krishna, "Flexagon: A multi-dataflow sparse-sparse matrix multiplication accelerator for efficient dnn processing," in *Proceedings of the 28th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 3*, ser. ASPLOS 2023. New York, NY, USA: Association for Computing Machinery, 2023, p. 252–265. [Online]. Available: <https://doi.org/10.1145/3582016.3582069>
- [8] M. Abolhasani and E. Kumacheva, "The rise of self-driving labs in chemical and materials sciences," *Nature Synthesis*, vol. 0, no. 0, pp. 1–10, 2023.
- [9] Q. Rao and J. Frtunikj, "Deep learning for self-driving cars: Chances and challenges," in *2018 IEEE/ACM 1st International Workshop on Software Engineering for AI in Autonomous Systems (SEFAIAS)*, ser. SEFAIS '18. New York, NY, USA: Association for Computing Machinery, 2018, p. 35–38. [Online]. Available: <https://doi.org/10.1145/3194085.3194087>
- [10] J. Jumper, R. Evans, A. Pritzel, T. Green, M. Figurnov, O. Ronneberger, K. Tunyasuvunakool, R. Bates, A. Žídek, A. Potapenko *et al.*, "Highly accurate protein structure prediction with alphafold," *Nature*, vol. 596, no. 7873, pp. 583–589, 2021.
- [11] X. Zhang, H. Ye, J. Wang, Y. Lin, J. Xiong, W. Hwu, and D. Chen, "DNNExplorer: A Framework for Modeling and Exploring a Novel Paradigm of FPGA-based DNN Accelerator," in *ICCAD*. New York, NY, USA: Association for Computing Machinery, 2020, pp. 1–9.
- [12] P. Xu, X. Zhang, C. Hao, Y. Zhao, Y. Zhang, Y. Wang, C. Li, Z. Guan, D. Chen, and Y. Lin, "Autodnnchip: An automated dnn chip predictor and builder for both fpgas and asics," in *Proceedings of the 2020 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*. New York, NY, USA: Association for Computing Machinery, 2020, p. 40–50. [Online]. Available: <https://doi.org/10.1145/3373087.3375306>
- [13] H. Ye, X. Zhang, Z. Huang, G. Chen, and D. Chen, "HybridDNN: A framework for high-performance hybrid DNN accelerator design and implementation," in *DAC*. San Francisco, CA, USA: IEEE, 2020, pp. 1–6.
- [14] H. Kwon, P. Chatarasi, V. Sarkar, T. Krishna, M. Pellauer, and A. Parashar, "Maestro: A data-centric approach to understand reuse, performance, and hardware cost of dnn mappings," *IEEE Micro*, vol. 40, no. 3, pp. 20–29, 2020.
- [15] Y.-H. Chen, T.-J. Yang, J. Emer, and V. Sze, "Eyeriss v2: A flexible accelerator for emerging deep neural networks on mobile devices," *IEEE Journal on Emerging and Selected Topics in Circuits and Systems*, vol. 9, no. 2, pp. 292–308, 2019.
- [16] TVM Developers, "VTA: Deep learning accelerator stack," 2020. [Online]. Available: [docs.tvm.ai/vta](https://docs.tvm.ai/vta)
- [17] J. Ngadiuba, V. Loncar, M. Pierini, S. Summers, G. Di Guglielmo, J. Duarte, P. Harris, D. Rankin, S. Jindariani, M. Liu *et al.*, "Compressing deep neural networks on FPGAs to binary and ternary precision with hls4ml," *ML: Science and Technology*, vol. 2, no. 1, p. 015001, 2020.
- [18] S. Skalicky, J. Monson, A. Schmidt, and M. French, "Hot & Spicy: Improving Productivity with Python and HLS for FPGAs," in *FCCM*. Boulder, CO, USA: IEEE, 2018, pp. 85–92.
- [19] N. Bohm Agostini, S. Dong, E. Karimi, M. T. Lapuerta, J. Cano, J. L. Abellán, and D. Kaeli, "Design space exploration of accelerators and end-to-end dnn evaluation with tflite-soc," in *SBAC-PAD*. Porto, Portugal: IEEE, 2020, pp. 10–19.
- [20] ARM Developers, "AMBA AXI and ACE Protocol Specification," 2020. [Online]. Available: <https://developer.arm.com/documentation/ih0022/e/AMBA-AXI3-and-AXI4-Protocol-Specification>
- [21] S. Liu, H. Fan, X. Niu, H.-c. Ng, Y. Chu, and W. LUK, "Optimizing cnn-based segmentation with deeply customized convolutional and deconvolutional architectures on fpga," *ACM Trans. Reconfigurable Technol. Syst.*, vol. 11, no. 3, dec 2018. [Online]. Available: <https://doi.org/10.1145/3242900>
- [22] C. Lattner, M. Amini, U. Bondhugula, A. Cohen, A. Davis, J. Pienaar, R. Riddle, T. Shpeisman, N. Vasilache, and O. Zinenko, "MLIR: Scaling Compiler Infrastructure for Domain Specific Computation," in *CGO*. Seoul, Korea (South): IEEE, 2021, pp. 2–14.
- [23] M. Developers, "'linalg' Dialect," 2020, online accessed on 11-04-2023. [Online]. Available: <https://mlir.llvm.org/docs/Dialects/Linalg/>
- [24] T. D. Le, G.-T. Bercea, T. Chen, A. E. Eichenberger, H. Imai, T. Jin, K. Kawachiya, Y. Negishi, and K. O'Brien, "Compiling ONNX Neural Network Models Using MLIR," *ArXiv*, vol. 0, no. 0, pp. 1–8, 2020.
- [25] Xilinx, "AXI Reference Guide," 2012. [Online]. Available: [https://docs.xilinx.com/v/u/14.1-English/ug761\\_axi\\_reference\\_guide](https://docs.xilinx.com/v/u/14.1-English/ug761_axi_reference_guide)
- [26] C. Salvador Rohwedder, N. Henderson, J. a. P. L. De Carvalho, Y. Chen, and J. N. Amaral, "To pack or not to pack: A generalized packing analysis and transformation," in *Proceedings of the 21st ACM/IEEE International Symposium on Code Generation and Optimization*, ser. CGO 2023. New York, NY, USA: Association for Computing Machinery, 2023, p. 14–27.
- [27] J. E. Stone, D. Gohara, and G. Shi, "Opencl: A parallel programming standard for heterogeneous computing systems," *Computing in Science and Engg.*, vol. 12, no. 3, p. 66–73, may 2010.
- [28] R. Reyes, G. Brown, R. Burns, and M. Wong, "Sycl 2020: More than meets the eye," in *Proceedings of the International Workshop on OpenCL*, ser. IWOCCL '20. New York, NY, USA: Association for Computing Machinery, 2020.
- [29] J. Haris, P. Gibson, J. Cano, N. Bohm Agostini, and D. Kaeli, "SECDATFLite: A toolkit for efficient development of FPGA-based DNN accelerators for edge inference," *Journal of Parallel and Distributed Computing*, vol. 173, pp. 140–151, 2023.
- [30] X. Jiao, Y. Yin, L. Shang, X. Jiang, X. Chen, L. Li, F. Wang, and Q. Liu, "Tinybert: Distilling bert for natural language understanding," *arXiv preprint arXiv:1909.10351*, vol. 0, no. 0, pp. 1–12, 2019.
- [31] J. Albericio, P. Judd, T. Hetherington, T. Aamodt, N. E. Jerger, and A. Moshovos, "Cnvlutin: Ineffective-neuron-free deep neural network computing," *ACM SIGARCH Computer Architecture News*, vol. 44, no. 3, pp. 1–13, 2016.
- [32] S. Zhang, Z. Du, L. Zhang, H. Lan, S. Liu, L. Li, Q. Guo, T. Chen, and Y. Chen, "Cambricon-x: An accelerator for sparse neural networks," in *2016 49th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*, IEEE. Taipei, Taiwan: IEEE, 2016, pp. 1–12.
- [33] Y. Chen, T. Luo, S. Liu, S. Zhang, L. He, J. Wang, L. Li, T. Chen, Z. Xu, N. Sun *et al.*, "Dadiannao: A machine-learning supercomputer," in *2014 47th Annual IEEE/ACM International Symposium on Microarchitecture*, IEEE. Cambridge, UK: IEEE, 2014, pp. 609–622.
- [34] C. Zhang, P. Li, G. Sun, Y. Guan, B. Xiao, and J. Cong, "Optimizing fpga-based accelerator design for deep convolutional neural networks," in *Proceedings of the 2015 ACM/SIGDA international symposium on field-programmable gate arrays*. New York, NY, USA: Association for Computing Machinery, 2015, pp. 161–170.
- [35] The Linux Perf Team, "Perf wiki," n.d., accessed on April 13, 2023. [Online]. Available: [https://perf.wiki.kernel.org/index.php/Main\\_Page](https://perf.wiki.kernel.org/index.php/Main_Page)
- [36] A. Developer, "Neon registers," 2023. [Online]. Available: <https://developer.arm.com/documentation/dht0002/a/Introducing-NEON/NEON-architecture-overview/NEON-registers>
- [37] F. Wang, M. Jiang, C. Qian, S. Yang, C. Li, H. Zhang, X. Wang, and X. Tang, "Residual attention network for image classification," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, Honolulu, HI, USA, 2017, pp. 3156–3164.
- [38] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz, J. Davison, S. Shleifer, P. von Platen, C. Ma, Y. Jernite, J. Plu, C. Xu, T. Le Scao, S. Gugger, M. Drame, Q. Lhoest, and A. Rush, "Transformers: State-of-the-art natural language processing," in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. Online: Association for Computational Linguistics, Oct. 2020, pp. 38–45. [Online]. Available: <https://aclanthology.org/2020.emnlp-demos.6>
- [39] Torch-MLIR Developers, "The Torch-MLIR Project," 2021. [Online]. Available: <https://github.com/llvm/torch-mlir>
- [40] E. Del Sozzo, R. Baghdadi, S. Amarasinghe, and M. D. Santambrogio, "A Unified Backend for Targeting FPGAs from DSLs," in *ASAP*. Milan, Italy: IEEE, 2018, pp. 1–8.
- [41] N. Bohm Agostini, S. Curzel, J. Zhang, A. Limaye, C. Tan, V. Amaty, M. Minutoli, V. G. Castellana, J. Manzano, D. Brooks, G.-Y. Wei, and A. Tumeo, "Bridging python to silicon: The soda toolchain," *IEEE Micro*, vol. 42, no. 5, 2022.

1026 [42] N. Bohm Agostini, S. Curzel, V. Amatya, C. Tan, M. Minutoli, V. G.  
1027 Castellana, J. Manzano, D. Kaeli, and A. Tumeo, "An mlir-based  
1028 compiler flow for system-level design and hardware acceleration," in  
1029 *ICCAD*. New York, NY, USA: Association for Computing Machinery,  
1030 2022.

1031 [43] A. Stjerngren, P. Gibson, and J. Cano, "Bifrost: End-to-End Evaluation  
1032 and optimization of Reconfigurable DNN Accelerators," in *2022 IEEE  
1033 International Symposium on Performance Analysis of Systems and  
1034 Software (ISPASS)*. Singapore: IEEE, May 2022, pp. 288–299.

1035 [44] S. Xiang, Y.-H. Lai, Y. Zhou, H. Chen, N. Zhang, D. Pal, and Z. Zhang,  
1036 "Heteroflow: An accelerator programming model with decoupled data  
1037 placement for software-defined fpgas," in *Proceedings of the 2022  
1038 ACM/SIGDA International Symposium on Field-Programmable Gate  
1039 Arrays*, ser. FPGA '22. New York, NY, USA: Association for  
1040 Computing Machinery, 2022, p. 78–88.

1041 [45] Y.-H. Lai, Y. Chi, Y. Hu, J. Wang, C. H. Yu, Y. Zhou, J. Cong,  
1042 and Z. Zhang, "Heterocl: A multi-paradigm programming infrastructure  
1043 for software-defined reconfigurable computing," in *Proceedings of the  
1044 2019 ACM/SIGDA International Symposium on Field-Programmable  
1045 Gate Arrays*, ser. FPGA '19. New York, NY, USA: Association for  
1046 Computing Machinery, 2019, p. 242–251.

1047 [46] X. Yang, M. Gao, Q. Liu, J. Setter, J. Pu, A. Nayak, S. Bell, K. Cao,  
1048 H. Ha, P. Raina, C. Kozyrakis, and M. Horowitz, "Interstellar: Using  
1049 halide's scheduling language to analyze dnn accelerators," in *Proceed-  
1050 ings of the Twenty-Fifth International Conference on Architectural Sup-  
1051 port for Programming Languages and Operating Systems*, ser. ASPLOS  
1052 '20. New York, NY, USA: Association for Computing Machinery,  
1053 2020, p. 369–383.

1054 [47] S. Dave, Y. Kim, S. Avancha, K. Lee, and A. Shrivastava, "Dmazerunner:  
1055 Executing perfectly nested loops on dataflow accelerators," *ACM Trans.  
1056 Embed. Comput. Syst.*, vol. 18, no. 5s, oct 2019.

1057 [48] J. Cong and J. Wang, "Polysa: Polyhedral-based systolic array  
1058 auto-compilation," in *2018 IEEE/ACM International Conference on  
1059 Computer-Aided Design (ICCAD)*. San Diego, CA, USA: IEEE, 2018,  
1060 pp. 1–8.

1061 [49] Q. Huang, M. Kang, G. Dinh, T. Norell, A. Kalaiah, J. Demmel,  
1062 J. Wawrzynek, and Y. S. Shao, "Cosa: Scheduling by constrained  
1063 optimization for spatial accelerators," in *2021 ACM/IEEE 48th Annual  
1064 International Symposium on Computer Architecture (ISCA)*. Valencia,  
1065 Spain: IEEE, 2021, pp. 554–566.

1066 [50] M. Vaidya, A. Sukumaran-Rajam, A. Rountev, and P. Sadayappan,  
1067 "Comprehensive accelerator-dataflow co-design optimization for convo-  
1068 lutional neural networks," in *2022 IEEE/ACM International Symposium  
1069 on Code Generation and Optimization (CGO)*. Seoul, South Korea:  
1070 Association for Computing Machinery, 2022, pp. 325–335.

1071 [51] S.-C. Kao, G. Jeong, and T. Krishna, "ConfuciuX: Autonomous hardware  
1072 resource assignment for dnn accelerators using reinforcement learning,"  
1073 in *2020 53rd Annual IEEE/ACM International Symposium on Microar-  
1074 chitecture (MICRO)*. Athens, Greece: IEEE, 2020, pp. 622–636.

1075 [52] S. Zheng, Y. Liang, S. Wang, R. Chen, and K. Sheng, "Flextensor:  
1076 An automatic schedule exploration and optimization framework for  
1077 tensor computation on heterogeneous system," in *Proceedings of the  
1078 Twenty-Fifth International Conference on Architectural Support for  
1079 Programming Languages and Operating Systems*, ser. ASPLOS '20.  
1080 New York, NY, USA: Association for Computing Machinery, 2020, p.  
1081 859–873.

1082 [53] T. Chen, L. Zheng, E. Yan, Z. Jiang, T. Moreau, L. Ceze,  
1083 C. Guestrin, and A. Krishnamurthy, "Learning to optimize tensor  
1084 programs," in *Advances in Neural Information Processing Systems*,  
1085 S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-  
1086 Bianchi, and R. Garnett, Eds., vol. 31. Curran Associates, Inc.,  
1087 2018. [Online]. Available: [https://proceedings.neurips.cc/paper\\_files/  
1088 paper/2018/file/8b5700012be65c9da25f49408d959ca0-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2018/file/8b5700012be65c9da25f49408d959ca0-Paper.pdf)

1089 [54] L. Zheng, C. Jia, M. Sun, Z. Wu, C. H. Yu, A. Haj-Ali, Y. Wang, J. Yang,  
1090 D. Zhuo, K. Sen *et al.*, "Ansor: Generating High-Performance tensor  
1091 programs for deep learning," in *14th USENIX symposium on operating  
1092 systems design and implementation (OSDI 20)*, 2020, pp. 863–879.

1093 [55] H. Min, J. Kwon, and B. Egger, "Flexer: Out-of-order scheduling  
1094 for multi-npus," in *Proceedings of the 21st ACM/IEEE International  
1095 Symposium on Code Generation and Optimization*, ser. CGO 2023.  
1096 New York, NY, USA: Association for Computing Machinery, 2023, p.  
1097 212–223.

1098 [56] M. Developers, "Transform Dialect: Fine-grain transformation control

dialect," 2022, online accessed on 11-04-2023. [Online]. Available: 1099  
<https://mlir.llvm.org/docs/Dialects/Transform/> 1100