

# MESC: Re-thinking Algorithmic Priority and/or Criticality Inversions for Heterogeneous MCSs

Jiapeng Guan<sup>†</sup>, Ran Wei<sup>‡†\*</sup>, Dean You<sup>§</sup>, Yingquan Wang<sup>†</sup>, Ruizhe Yang<sup>†</sup>, Hui Wang<sup>§</sup>, Zhe Jiang<sup>§\*</sup>  
<sup>†</sup>Dalian University of Technology, China. <sup>‡</sup>Lancaster University, UK. <sup>§</sup>Southeast University, China.

**Abstract**—Modern Mixed-Criticality Systems (MCSs) rely on hardware heterogeneity to satisfy ever-increasing computational demands. However, most of the heterogeneous co-processors are designed to achieve high throughput, with their micro-architectures executing the workloads in a streaming manner. This streaming execution is often non-preemptive or limited-preemptive, preventing tasks’ prioritisation based on their importance and resulting in frequent occurrences of algorithmic priority and/or criticality inversions. Such problems present a significant barrier to guaranteeing the systems’ real-time predictability, especially when co-processors dominate the execution of the workloads (e.g., DNNs and transformers).

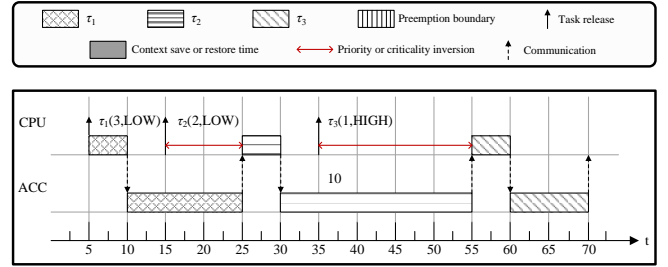
In contrast to existing works that typically enable coarse-grained context switch by splitting the workloads/algorithms, we demonstrate a method that provides fine-grained context switch on a widely used open-source DNN accelerator by enabling instruction-level preemption without any workloads/algorithms modifications. As a systematic solution, we build a real system, i.e., Make Each Switch Count (MESC), from the SoC and ISA to the OS kernel. A theoretical model and analysis are also provided for timing guarantees. Experimental results reveal that, compared to conventional MCSs using non-preemptive DNN accelerators, MESC achieved a 250x and 300x speedup in resolving algorithmic priority and criticality inversions, with less than 5% overhead. To our knowledge, this is the first work investigating algorithmic priority and criticality inversions for MCSs at the instruction level.

## I. INTRODUCTION

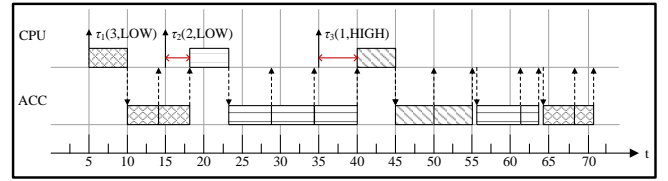
Mixed-Criticality Systems (MCSs) are safety-critical systems in which functionalities are developed at multiple assurance/criticality levels and integrated on a shared platform, e.g., System-on-Chip (SoC) [1]–[5]. For instance, in the automotive industry, an Advanced Driver Assistance System (ADAS) may involve various functionalities developed at different criticality levels, with collision avoidance being of high criticality, whereas route-planning may be of lower criticality [6]–[8].

To satisfy the computational demands of diverse functionalities in MCSs, e.g., running workloads on Deep Neural Networks (DNNs) or transformers, hardware vendors have developed SoCs with different architectures [9]–[11], with a high degree of *heterogeneity*. That is, the SoCs couple general-purpose CPUs with heterogeneous co-processors to accelerate algorithmic executions [12]–[16].

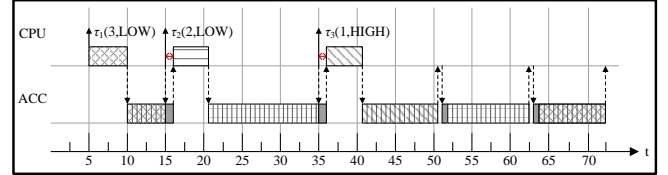
**Research challenges.** The deployment of these co-processors brings frequent occurrences of *algorithmic priority and/or*



(a) Scheduling with a non-preemptive DNN accelerator.



(b) Scheduling with a DNN accelerator with limited preemption [17]–[21].



(c) Scheduling with DNN accelerator allowing instruction-level preemption.

Fig. 1. Scheduling with DNN accelerators (referred to as ACC in figures) featuring different preemption characteristics: (a) non-preemption, (b) limited preemption, and (c) instruction-level preemption.  $\tau_i(P_i, L_i)$  represents task  $\tau_i$  with priority  $P_i$  and criticality  $L_i$ , with smaller  $P_i$  indicating a higher priority.

*criticality inversions*<sup>1</sup>, posing a significant challenge to guaranteeing system-wide real-time performance [17]–[20], [29], [30]. Specifically, unlike CPU micro-architecture, which provides fine-grained<sup>2</sup> context switches, the micro-architecture of the co-processors is typically designed to maximise throughput while running designated parallel and computation-intensive workloads [12]–[16], [32]. This often results in the computation being processed in a streaming order, which is either non-preemptive

<sup>1</sup>Due to the differing perspectives on the importance of low-criticality tasks between academia and industry, practical systems cannot tolerate the abandonment of these tasks in degraded modes [22]–[28], which may cause criticality inversion. More detailed explanations will be provided in Sec. II.

<sup>2</sup>The ability of the OS kernel to trigger a context switch after the commitment of every single instruction, if no critical section is being executed [31].

\* represents corresponding authors. Emails: r.wei5@lancaster.ac.uk and zhejiang.uk@gmail.com.

or only preemptive at the boundaries of the algorithm<sup>3</sup> [17]–[21] (i.e., limited preemption). Such hardware limitations bar the Operating System (OS) from creating appropriate control flow for context switches in co-processors [33], [34]. Hence, without modifications to the co-processor hardware and/or workloads/algorithms to enhance preemption efficiency, the prolonged occupation of shared resources hinders the effective application of traditional scheduling methods, such as fixed-priority [35], [36] and Earliest-Deadline-First [37]–[39], in ensuring schedulability by prioritising task execution based on task importance [40]–[42]. That is, whilst a less important task occupies the co-processor, more important tasks must wait until the end of the co-processor’s ongoing computation.

Figs. 1(a) and 1(b)<sup>4</sup> provide an example through a heterogeneous MCS with a CPU and a DNN accelerator. When the accelerator is non-preemptive (Fig. 1(a)), tasks experience significant priority and/or criticality inversions [47]–[49], as high-priority/criticality tasks have to be postponed until the accelerator completes all ongoing computations. When the accelerator is designed/configured to allow limited preemption using existing techniques [17]–[20], [29], [30] (Fig. 1(b)), the capacity of context switches becomes coarse-grained. However, since preemption is only possible at the boundaries of the executed algorithm, the system still suffers from considerable priority and/or criticality inversions.

**Contributions.** Here, we introduce **Make Each Switch Count (MESc)**, a heterogeneous MCS framework that minimises algorithmic priority and criticality inversions down to the instruction level (Fig. 1(c)). To achieve this, we present

- a new DNN accelerator (Gemmini<sup>RT</sup>) based on the open-source NPU<sup>5</sup> architecture [15], [16], which enables instruction-level preemption. This lays the foundation for fine-grained context switches in heterogeneous MCSs;
- a context switch strategy and OS add-ons to ensure data consistency and manage DNN accelerator context switches. This provides software-level support for DNN accelerator context switches;
- a full-stack framework that incorporates customised SoC and Instruction Set Architecture (ISA) to OS kernel, forming a complete solution for heterogeneous MCSs;
- a theoretical model and analysis for the proposed framework provide theoretical validation for MESc.

We deployed our proposed system on AMD Alveo U280 FPGA and examined it using various metrics, including blocking duration, real-time performance, and overhead. Experiments show that compared to conventional non-preemptive DNN accelerators, Gemmini<sup>RT</sup> achieved 250x and 300x accelerations in resolving algorithmic priority inversion and criticality inversion, respectively. Furthermore, deploying the Gemmini<sup>RT</sup>

<sup>3</sup>The boundary of an algorithm is typically defined as the completion of an operator (e.g., Softmax and ReLu) or a part of the input data.

<sup>4</sup>Fig. 1 is a simplified example. In practice, co-scheduling between CPUs and co-processors involves a wide range of strategies [43]–[46], yet most of them still encounter challenges related to priority and/or criticality inversions.

<sup>5</sup>A specialised processor designed to accelerate machine learning and deep learning tasks.

in a heterogeneous SoC can significantly improve the timing performance of MCSs with negligible hardware overhead.

## II. PRELIMINARIES

### A. Dual-mode MCS

Conventional MCS theoretical models often assume that the Worst-Case Execution Time (WCET) of a task is estimated with varying degrees of confidence [9], [50], [51]. A task’s high-critical WCET (HI-WCET) is associated with a high level of confidence but tends to be overly pessimistic, whereas its low-critical WCET (LO-WCET) is less pessimistic but has a lower confidence level. The correctness criteria specify that if all tasks finish execution within their LO-WCETs, then they will all finish execution by their deadlines. However, if any high-critical task’s (HI-task’s) execution time exceeds its LO-WCET, the HI-tasks must complete execution by their deadlines [50], [52], [53]. To satisfy the above criteria, *mode switch* is a straightforward strategy: initially, the system operates in low-critical mode (LO-mode), in which the system operates under the assumption that the execution time for all tasks will not exceed their LO-WCETs. If this assumption is violated (i.e. if any task fails to finish before its LO-WCET), the system *switches* to high-critical mode (HI-mode). In HI-mode, the scheduling policy may permit HI-tasks to execute beyond their LO-WCETs but will ensure that their execution times do not exceed their HI-WCETs. Consequently, to ensure that HI-tasks meet their deadlines, low-critical tasks (LO-tasks) may need to be terminated or executed with minimal time budget. The direct termination of LO-tasks could lead to potential safety hazards, as the system may still rely on LO-tasks to perform as expected when entering HI-mode [22]–[28]. Thus, similar to Imprecise MCSs [34], [54], we do not drop LO-tasks, instead, we continue to execute LO-tasks in HI-mode, but only when none of the HI-tasks are being executed.

### B. The Gemmini NPU Architecture

Gemmini [16] is part of the open-source RISC-V ecosystem [55], developed for machine learning by UC Berkeley. The key architecture is based on a systolic array, including multiple tiles. Each tile consists of a configurable number of Processing Elements (PEs). The array reads data from a local, explicitly managed scratchpad of banked SRAMs, and writes results or intermediate values to a local accumulator.

**Gemmini instructions.** Gemmini is designed to operate in conjunction with a RISC-V CPU core by executing bespoke, well-defined instructions. These instructions are dispatched by the CPU and initially received by Gemmini’s central controller, forwarded to a reservation station for classification. The reservation station categorises the instructions into configuration, load, store and compute. Configuration instructions are directly executed, while others are issued to their respective controllers. These controllers decompose the instructions and delegate them to different functional modules for execution.

**Other co-processors.** As discussed in Sec. I, algorithmic priority and criticality inversions are common problems for the MCSs built on heterogeneous SoCs. Here, we choose Gemmini as the co-processor as our case study for two reasons: (i)

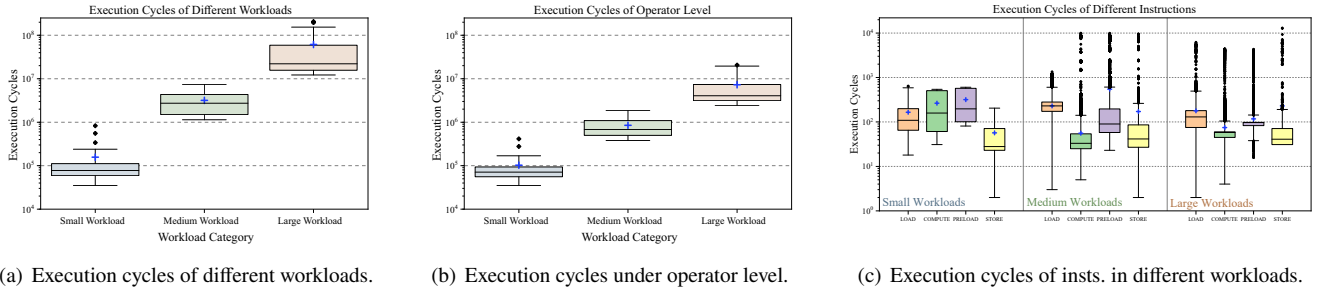


Fig. 2. Execution cycles of Gemmini running workloads of varying sizes. The workloads are categorised based on their execution times. Small workload: [0, 1 million] cycles; medium workloads: (1 million, 10 million] cycles; large workloads: (10 million, 1 billion] cycles.

from the micro-architecture perspective, Gemmini features a representative design similar to other NPUs [56]–[58], facing the same challenges we identified in this work; (ii) from the SoCs/systems perspective, Gemmini is a popular DNN accelerator that has been integrated into various SoCs [59]–[62] and supports modern machine-learning algorithms (e.g., DNNs and transformer), which is crucial for the community.

### III. MOTIVATIONS AND CHALLENGES

In this section, we demonstrate the need for fine-grained context switch in heterogeneous MCSs through the quantitative illustration of the penalty due to algorithmic priority and/or criticality inversions. We first examine the blocking duration caused by DNN accelerators with different preemptive capacities, then discuss our rationale for selecting MCS as the subject of our studies, followed by an exploration of the challenges in building MESec.

#### A. Quantitative Analysis of Priority and Criticality Inversions

To quantitatively assess the blocking duration of priority and criticality inversions potentially introduced by DNN accelerators with different preemptive capacities, we execute a collection of workloads of different sizes, using DNNs including AlexNet, MobileNet, ResNet50 and Transformer. We developed a cycle-accurate performance counter and integrated it into the instruction commit stage, driven by the same clock source as Gemmini. This allows us to measure and record the cycles for entire workloads, individual operators (i.e., independent and executable computation blocks or linear operators, such as Softmax and ReLU, within the workload), and single instructions (as illustrated in Fig. 2), thus reflecting the impact of varying preemption capabilities of different DNN accelerators on the blocking duration of high-priority/criticality tasks.

**Non-preemption.** When the DNN accelerator is non-preemptive, tasks may endure algorithmic priority and/or criticality inversions, with durations that span from tens of thousands of cycles, to several hundred million cycles, as shown in Fig. 2(a). In safety-critical applications, such prolonged “obligatory waiting periods” may cause delays to critical

functions that lead to hazardous events, consequently causing accidents that cause harm.

**Limited preemption.** When the DNN accelerator is featured with limited preemptive capacity, the “obligatory waiting periods” can be shortened. However, the remaining execution time of several million cycles still poses a significant challenge for practical applications, as shown in Fig. 2(b).

**Instruction-level preemption.** When the DNN accelerator enables preemption at the instruction level (as a result of this work, detailed in Sec. VIII), the longest execution time for instructions is 2 orders of magnitude shorter than the “obligatory waiting periods” in DNN accelerator with limited context switch (even in the worst-case scenario). As shown in Fig. 2(c), we classified the instructions into four types<sup>6</sup> and plotted their execution cycles.

Hence, with the provision of instruction-level switches, the duration of priority and criticality inversions could be reduced to the time needed for instruction and context switch, allowing timely preemption for high-priority/criticality tasks.

#### B. MCS: an Example

Although contributions to priority inversion issues are applicable to a wide range of real-time systems, we have chosen MCS as our system model for the following reasons.

**System complexity.** In the context of algorithmic priority inversion, MCS introduces an additional challenge: criticality inversion [22]–[28], [34], [54], [63]. In this regard, MCS is not only a typical complex system but also a representative of systems with broader real-time system challenges (i.e., other real-time systems can be viewed as subsets of MCS). By addressing both the broader issue of priority inversion in real-time systems and the more specific problem of criticality inversion in MCS, we demonstrate the effectiveness of our approach, without limiting its applicability to other systems.

**Popularity.** MCS has remained a vital and popular research area for two decades, attracting continuous contributions and

<sup>6</sup>The preload instruction [15], [16] is used to prefetch data into the internal memory (e.g., scratchpad and accumulator), working in conjunction with compute instructions to optimise data transfer and pipeline execution.

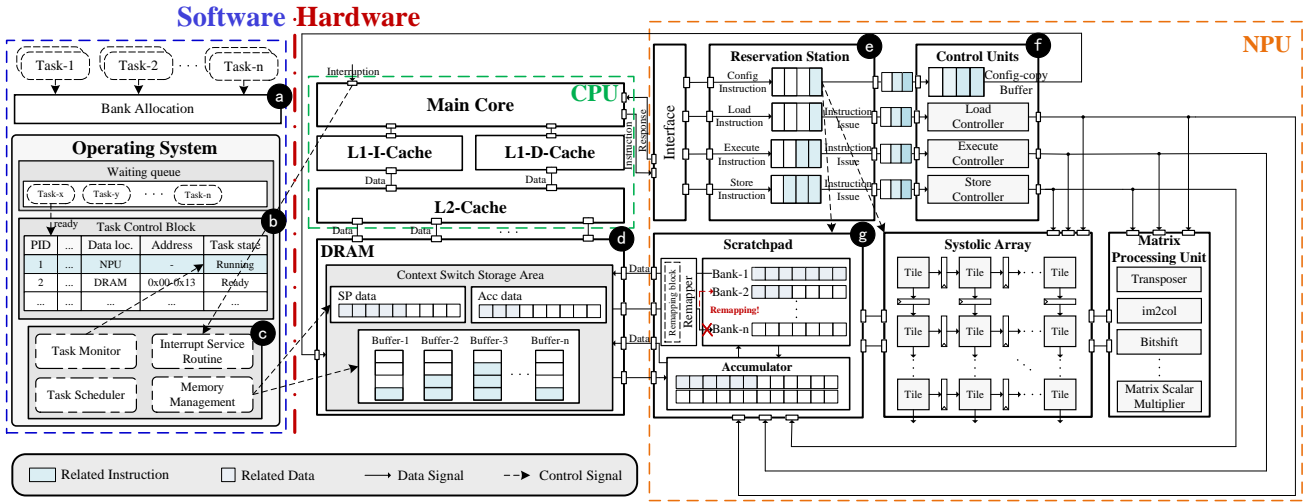


Fig. 3. Architectural overview: the blue box represents the software level, the orange and green boxes correspond to the NPU and CPU (hardware), respectively.

advancements from the research community [2], [7], [25], [34], [53], [64]–[67], and playing a crucial role in industries with high safety-critical requirements [1]–[5] (as detailed in Sec. I).

### C. Research Challenges

Based on our observations, it is crucial to enable instruction-level context switches for DNN accelerators to minimise algorithmic priority and criticality inversions. However, achieving this is not straightforward, as several key challenges must be addressed:

- **Complexity.** It is important to overcome the inherent complexity involved in DNN accelerator context switch, to properly handle the saving and loading of computation and configuration data to achieve comprehensive accelerator context switches (and develop an acceleration strategy);
- **Scheduling.** The OS scheduler shall be enhanced to orchestrate context switch in the accelerator to correctly schedule tasks and provide software-level support for context switch so that the accelerator switch becomes transparent.
- **Hardware and software integration.** System-level integration of hardware and software is needed for a generic approach in building heterogeneous MCSs with DNN accelerators that support instruction-level context switches.
- **Modelling and analysis.** Theoretical modelling of the system is needed for a comprehensive theoretical analysis to show that our proposed approach is fit for purpose.

## IV. MESC: THE SYSTEMATIC FRAMEWORK

In coping with the above challenges, we have developed MESC, which includes a DNN accelerator (called Gemmini<sup>RT</sup>), allowing instruction-level preemption, and OS add-ons that provide software-level support for context switch control flow. The architectural overview of MESC is given in Fig. 3.

**Hardware level.** At the hardware level, we have developed a context switch and acceleration mechanism for the DNN accelerator. During the context-saving process, computation data is directly transferred from the scratchpad to the DRAM

(Fig. 3. **d**), and the address information is recorded in the task’s Task Control Block (TCB, Fig. 3. **b**). For configuration data, we employ a config-copy buffer (Fig. 3. **f**) to store the most recently executed configuration instructions of different types, and these are also sent to the DRAM (Fig. 3. **d**). During the context-restoring process, computation data is re-loaded into the local memory of the accelerator (Fig. 3. **g**) using addresses recorded in the TCB. The accelerator will also be reconfigured by the stored configuration instructions. Instructions previously dispatched by the CPU without receiving a response signal from the accelerator are then resent to it. Also, we implemented an address remapper (Fig. 3. **g**), which centralises the storage of data into allocated banks and uses a remapping block to maintain this process, thereby supporting the subsequent local memory<sup>7</sup> allocation methods we employ. The remapping block is sent back to the DRAM during each context switch and updated during each context restoration, based on the bank addresses into which the computation data is reloaded. When there is sufficient local memory, a minimal amount of computation data needs to be preserved, reducing the time required for context switches. Lastly, we have added specific ISA to the system, (Tbl. I), which we will use later in Secs. V and VI.

**Software level.** At the software level, we provide an accelerator local memory allocation method (Fig. 3. **a**), along with a task monitor and scheduler (Fig. 3. **c**) that can be adopted by (real-time) OS kernels. This method, which is enhanced by the address remapper (Fig. 3. **g**), adjusts the local memory allocated to tasks, thereby directly implementing hardware segmentation to reduce the time required for context switches when local memory is sufficient. This is under the fundamental premise of maintaining unchanged task execution times. We also added the task monitor and scheduler to the OS kernel space. The monitor uses hardware timers to track the running conditions of tasks, while the scheduler provides control functions, e.g.,

<sup>7</sup>The “local memory” allocated by this method specifically refers to the scratchpad in this paper.

TABLE I  
NEW ISA ADDED FOR GEMMINI<sup>RT</sup>.

Name/Type	Description
<b>instruction_freeze</b>	Halts the execution of all queued insts., except for flush-related insts.
<b>step_wise_mv</b> <b>in</b> <b>/mvout</b>	Moves data into/out of the scratchpad using the configuration channel.
<b>mvin/mvout_config_buffer</b>	Moves stored configuration insts. from/to DRAMs.
<b>reconfig</b>	Executes and clears insts. in the config-copy buffer.
<b>mvin/mvout_remapping_block</b>	Moves the remapping block into/out of the accelerator using the default configuration channel.
<b>flush_x</b>	The x represents freeze, bank, etc., either resuming the execution of insts. or flushing data within a specified component.

Context\_switch and Mode\_switch to manage tasks and mode switches. In the OS user space, we preserve compatibility with existing OS interfaces, enabling the migration and adaptation of tasks initially designed for traditional MCS.

**Context switch.** A context switch is often triggered by an interrupt or a system call, where the OS transitions to the Interrupt Service Routine (ISR) to manage the current situation. If the interrupt is triggered by a timer monitoring the task (other interrupts are handled by the original OS kernel), control reverts to the task scheduler (Fig. 3. **c**). The scheduler then determines whether a context switch is necessary. If a context switch is required (detailed in Sec. VI), it initially prohibits the execution of instructions that are queued in the accelerator except flush instructions and waits for the completion of multiple instructions that are concurrently executing (Fig. 3. **a** and **f**). The scheduler then flushes the accelerator instruction queues (Fig. 3. **e** and **f**), ensuring no instructions are present in the accelerator except config-copy buffer. It reinstates the issuing and operation of instructions in the accelerator to ensure that the instructions related to context switch are not impeded. It then directs the accelerator to save the current context to DRAM (Fig. 3. **d**), including the computation and configuration data from the accelerator. When a task completes and a previously preempted task needs to resume execution, a similar approach is adopted. The relevant data is retrieved from DRAM and restored to the accelerator, the remapping block is updated, and configuration of the accelerator is returned to its previous state. Subsequently, the OS directs the CPU to re-dispatch the instructions that were previously dispatched to the accelerator but did not receive a response signal. This step restores the instructions that were queued in the accelerator. Following this procedure, the accelerator is fully restored to its state prior to the context switch.

**Mode switch.** The strategy of mode switch in MESC is governed by the following rules:

- **LO-mode:** the system prioritises the highest-priority tasks for execution. Under the bank allocation method, all tasks are executed using minimal accelerator local memory necessary to maintain their operational speed.
- **Mode transition:** when a HI-task exceeds its LO-WCET, the system commences mode transition. The system prioritises the scheduling of HI-tasks, and LO-tasks are scheduled only

when all HI-tasks are idle. During transition, only LO-tasks with computation data not yet saved back to the main memory are permitted to execute until the accelerator local memory contains data from at most one LO-task.

- **HI-mode:** once the accelerator local memory contains data from at most one LO-task, the system transitions to HI-mode. The system prioritises the scheduling of HI-tasks with accelerator local memory allocated according to the bank allocation method. When LO-tasks preempt each other, all related data must be evacuated from accelerator to ensure that the accelerator local memory contains data from at most one LO-task at any time, thereby reducing the average criticality inversion duration for HI-tasks. Moreover, when there are no tasks currently executing within the system, it will revert to LO-mode.

**Framework availability.** The MESC system architecture is designed with a high degree of availability. Specifically, (i) the architecture can be adapted with minor modifications to other DNN accelerators (e.g., Tetris [56], Sigma [57], and DNPU [58]); (ii) modifications to the OS retain the Application Programming Interfaces (APIs) used in traditional MCS frameworks, allowing user programs developed for traditional MCS frameworks to be directly ported to this system without major changes; (iii) in MESC, the CPU primarily functions as a instruction dispatcher to the accelerator. However, our approach can still be applied to systems employing other co-scheduling strategies to address challenges arising from priority and criticality inversions; (iv) although our research focuses on MCS, it is general enough to be applicable to other real-time systems as well; and (v) the MESC can be ported to other systems without altering the programming model. Specifically, the Gemmini programming model, which includes a high-level model [15], [16] for reading DNN descriptions from the Open Neural Network Exchange (ONNX) format and generating software binaries, as well as a low-level model [15], [16] for invoking library functions provided by UC Berkeley, remains unchanged.

## V. GEMMINI<sup>RT</sup>: THE MICRO-ARCHITECTURE

As introduced, we use Gemmini as a case study to demonstrate the efficacy of our framework, leading to the evolution of Gemmini into what we call Gemmini<sup>RT</sup> with the integration of a context switch mechanism. To do so, we outline three key requirements: Gemmini<sup>RT</sup> should (i) maintain the consistency of computation data, (ii) maintain the consistency of configuration data, and (iii) allocate hardware resources efficiently and optimise the context switches' overhead.

In response to these requirements, we designed a novel micro-architecture that includes (i) a default configuration channel that operates independently of the current accelerator configuration environment; (ii) a config-copy buffer that stores the most recent types of configuration instructions; and (iii) an address remapper that assists the accelerator local memory allocation.

### A. Handling Computation Data

In Gemmini, the matrix processing units (e.g., transposer, bit shifter) function as operational modules that terminate upon completion of their instructions and do not store data. Similarly, the systolic array can also be seen as a specialised computational unit. It executes operations based on instructions, storing intermediate results in the accumulator upon completing a single computation [15], [16]. Subsequent instructions retrieve values from the accumulator rather than relying on registers of the systolic array. Thus, all computation data is stored exclusively in the scratchpad and the accumulator. Data and memory management between the accelerator and the host CPU is explicit in Gemmini, meaning that data must be explicitly moved between the main address space of the processor and the private address space of the accelerator using a series of move instructions. In the original ISA, two data movement instructions, `mvin` and `mvout`, facilitate the transfer of data between the main memory and the internal accelerator via the Direct Memory Access (DMA) unit.

However, it is important to note that configuration instructions may influence these move instructions. If `mvin` and `mvout` are executed on an actively running accelerator, the data saved and restored could be incomplete. Therefore, to ensure configuration consistency and to secure the complete retrieval of computation data, we have introduced new move instructions, `step_wise_mvin` and `step_wise_mvout`, as shown on the right of Fig. 4. These new instructions utilise the default configuration channel (shown on the left of Fig. 4), which we developed within the control units of the accelerator, enabling the complete saving and restoring of computation data without altering the current accelerator configuration environment. Specifically, during the context-saving process, by utilising these move instructions and recording the target main memory addresses before the data transfer, we can ensure the complete preservation of computation data in the main memory through the scheduling of the OS. Similarly, during the context-restoring process, by retrieving the main memory address index of the computation data, we can fully load it back into the accelerator. Through such processes, we can completely save the computation data, and restore it as needed at any time.

### B. Handling Configuration Data

In Gemmini, configuration instructions are unique compared to other types of instructions, for a couple of reasons. Firstly, configurations are categorised into four types: load, store, execute, and norm, where a subsequent configuration instruction of the same category can override the previous one. Secondly, once a configuration instruction is dispatched by the CPU to the accelerator and enters the reservation station, it does not need to be issued to the execution units. Instead, it is executed immediately within the reservation station, directly affecting the corresponding registers to achieve the configuration purpose, with a 2-cycle execution time.

To accommodate the characteristics outlined above, we have implemented a dedicated buffer called the config-copy buffer to store the four most recent configuration instructions of

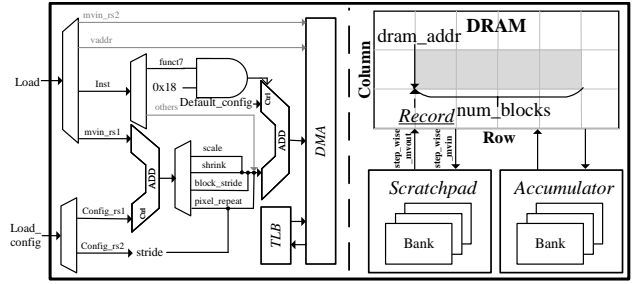


Fig. 4. The default configuration pathway for load class instructions (shown on the left, in black), and the process utilising `step_wise_mvin` and `step_wise_mvout` to transfer matrices into or out of the scratchpad and accumulator (right side).

each type. During the context-saving process, these instructions are transmitted back to DRAM using the newly added `mvout_config_buffer` instruction, with their addresses recorded by the monitor. During the context-restoring process, these instructions are re-loaded into the accelerator using the `mvin_config_buffer` instruction, and executed after all data has been fully reloaded to reconfigure the accelerator.

In scenarios involving complex system configurations, it is feasible to implement and maintain multiple high-level configuration instructions to streamline system setups and manage configurations. If such strategies prove insufficient due to overly intricate configurations, an alternative viable method is to introduce a flag bit for cache lines in the L1 I-cache of the CPU. Configuration instructions executed during task operations would update these flag bits, and the scheduler would track their locations, ensuring these cache lines are not displaced during cache entry switching. When tasks are finished or terminated, the associated cache line flags are cleared. During context restoration, these flagged cache lines are prioritised for execution. However, as noted in our Gemmini case study, due to the standardised nature and minimal quantity of configuration instructions, we opted not to employ this method.

### C. The Address Remapper

Given that not every task requires all accelerator resources, to enable more efficient resource utilisation, we introduced an address remapper to transition the explicit resource allocation in Gemmini to a semi-explicit form. This allows for effective hardware partitioning when local memory is sufficient, eliminating the need to save and restore scratchpad data during context switches and accelerating the context-switching process.

The address remapper, as shown in Fig. 5, operates by intercepting the DMA streams entering the scratchpad, incorporating a dynamic offset to relocate them to the appropriate bank (Fig. 5. **b** and **d**). This ensures that, externally, it appears as though the data is still being directed to its initial position within the scratchpad. We established a configurable-sized remapping block, defaulting to 4KB, within the scratchpad (Fig. 5. **c**) that stores the mapping from their original addresses to actual addresses, facilitating address translation and data management.



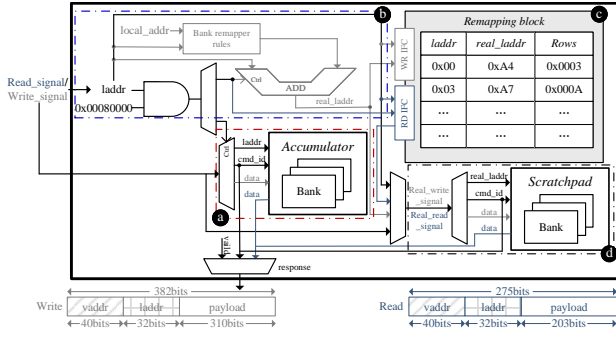


Fig. 5. Micro-architecture of the address remapper, featuring the signal `laddr` for the local memory address to be written or read, and `local_addr` representing the address calculation rules within the accelerator. Blue lines: the path of the read signals; grey lines: the path of the write signals; black lines: the common pathways.

Specifically, we have assigned a semaphore to each bank, referred to as the banklock, which is engaged when the bank is loaded with valid data. In the event of a task anomaly or task completion, during the context-switching process, the OS identifies which bank is involved by consulting the address mappings in the remapping block. It then deactivates the corresponding semaphore and flushes the data contained within that bank. When the DMA stream initiates a write request to the scratchpad, the system intercepts the stream and applies a dynamic offset to the scratchpad address bits before the stream enters the write queue (Fig. 5. **b**). This offset positions the data into a scratchpad bank that is either partially filled and locked by the task, or into one that is currently unlocked. The system then locks the bank and records the start and end positions of the address in the remapping block (Fig. 5. **c**). It is crucial, however, to ensure that the number of scratchpad banks that are accessed does not exceed the total number allocated by the system. Similarly, when the DMA stream requests a read from the scratchpad, the system first consults the remapping block (Fig. 5. **c**) to locate the corresponding mapping and then adjusts the scratchpad address bits (Fig. 5. **d**) in the DMA stream. When a task is preempted and subsequently resumed, the remapping block will adjust the target address based on the newly allocated bank that the task enters. In fact, in Gemini<sup>RT</sup>, the accumulator is almost identically structured to the scratchpad. However, compared to the scratchpad, the memory capacity of the accumulator is significantly smaller. Based on this, we have chosen not to impose allocation restrictions on the accelerator’s accumulator (Fig. 5. **a**).

Under this mechanism, the OS only needs to keep track of the number of banks currently allocated to schedule tasks effectively. The specific banks and addresses allocated are determined automatically by the hardware, ensuring efficient resource management and task execution.

## VI. THE OS KERNEL ADD-ONS

To ensure compatibility with the hardware mechanisms for context switches, we developed two software add-ons for real-time OS kernels (e.g., FreeRTOS) running on an open-source

### Algorithm 1: Function for Context Switch

```

/* Kernel.Scheduler:Context Switch */
1 Function Context_switch(task *current):
2   task *next = NULL
3   Kernel.Intr.Disable()
4   next = Kernel.Scheduler.Find_next_task()
5   if (current == NULL) then
6     continue;
7   else
8     if (current == next) then
9       continue
10    else
11      ACC.Instructions_freeze()
12      while (Kernel.Instructions_complete() ==
13        False) do
14        continue
15      end
16      Monitor.Timer.Pause(current)
17      Kernel.Scheduler.Context_save (current, next)
18    end
19    if (next != NULL and next != current) then
20      if (Monitor.Timer.Is_zero(next)) then
21        Monitor.Timer.Set(next)
22      else
23        Kernel.Scheduler.Context_restore (current,
24          next)
25      end
26      Monitor.Timer.Activate(next)
27      current = next
28    end
29    Kernel.Intr.Enable()
30    Kernel.Scheduler.Context_jump_to_PC(current)
31 End
/* Kernel.Scheduler:Context Save */
32 Function Context_save(task *current, task *next):
33   ACC.Preprocess_handler()
34   ACC.Accumulator_step_wise_mvout()
35   ACC.Mvout_config_buffer()
36   if (next->Bank + Kernel.Scheduler.Locked_banks() ≤
37     Kernel.Scheduler.Total_banks()) then
38     ACC.Banks_step_wise_mvout()
39     ACC.Release_banks(current)
40   end
41   ACC.Flush()
42 End

```

RISC-V processor (e.g., Rocket core [68]): (i) a task monitor to oversee task properties (including data locations) and key system checkpoints, and (ii) a task scheduler to help manage mode switches and context switches.

#### A. The Task Monitor

When an interrupt occurs, the ISR calls the task monitor to add or update task information within the TCB, storing task information, i.e., basic attributes, a program counter, data storage locations (DRAM or accelerator local memory) along with their addresses, and their status. The storage location of the task data determines if the computation data needs to be reloaded into the scratchpad during context restoration. A task’s current status can be: ready, running, pending, or interrupted.

Similar to traditional dual-criticality MCS frameworks, the task monitor creates a timer for each task. Specifically, upon system initialisation, the LO-WCETs of tasks are preloaded into memory. During context switches, the monitor suspends the current task’s timer and activates the next task’s timer. If a HI-task exceeds its LO-WCET in LO-mode, the timer will trigger an interrupt to execute the `Mode_switch` function to transition the system mode. Upon task completion, this timer is reset.

### B. The Task Scheduler

The task scheduler follows the system rules of MESC, as outlined in Sec. IV. Its responsibilities include: (i) the scheduling of an idle task when no other tasks are present, continuously checking for tasks that require execution; (ii) the routine maintenance of system status and assessing the necessity for context switches; (iii) the invocation of the `Context_switch` function and its auxiliary functions to assist the accelerator in performing context switches; and (iv) upon completion of the current task, flushing its data in the accelerator and identification of the next task that requires execution.

During task execution, the scheduler periodically checks for arriving tasks and maintains the system status using a hardware timer. The pseudocode in Alg. 1 details the steps involved for context switches. Initially, the kernel disables interrupts and the scheduler identifies the next task to execute. The scheduler then evaluates how to manage the context based on the current system task status. Specifically, if no tasks are executing, there is no need to save the context (Alg. 1: lines 5 - 6). If a task is executing, the scheduler checks whether it is the same as the current task. If they are the same, no further action is taken and interrupts are re-enabled before returning (Alg. 1: lines 8 - 9, 28 - 30). In other scenarios where accelerator context needs to be saved, the scheduler first prohibits the execution of instructions in the accelerator except for `flush` instructions which are neither issued in the reservation station nor started in the control module (Alg. 1: line 11). The kernel waits for response signals to ensure that all executing instructions are complete (Alg. 1: lines 12 - 14). Subsequently, the task monitor pauses the timer of the current task. The scheduler invokes the `Context_save` function to save the current context.

In `Context_save`, preprocessing is initially performed. Specifically, all instruction queues within the accelerator, except for the config-copy buffer, are flushed, and the issuing and execution of instructions are resumed. Computation data and configuration instructions are then preserved (Alg. 1: lines 33 - 34). The handling of data within the scratchpad is then determined by whether the number of banks in the scratchpad is sufficient to support the normal operation of subsequent tasks (Alg. 1: lines 35 - 38). Finally, the accelerator flushes all related data, including instructions in the queue, Translation Lookaside Buffer (TLB), and configuration data.

After the current context has been processed, if a task needs execution and data restoration, the `Context_restore` function is invoked (Alg. 1: lines 19 - 24). `Context_restore`’s logic mirrors `Context_save`. It decides whether to re-load task data based on the stored addresses, updates the remapping

block, re-configures the accelerator as necessary, and re-sends instructions that did not receive a response signal previously. Following context restoration, the monitor activates the timer for the task. Finally, the kernel enables interrupts and resumes task execution. It is important to note that this discussion only addresses scenarios that require a context switch. Situations that do not lead to a context switch are managed in the upper-level function of `Context_switch` within the scheduler.

Despite the introduction of a new system architecture in MESC, the design of the context switch minimises modifications, which are limited to slight adjustments to the TCB and the addition of a task monitor and a task scheduler, to the OS. Additionally, this design maintains the original OS APIs used in traditional MCS frameworks.

## VII. TIMING ANALYSIS AND OPTIMISATION

In this section, we present the Worst-Case Response Time (WCRT) analysis for MESC to assess the schedulability of a given task set with constrained deadlines. We also propose a local memory allocation method for distributing memory resources within the accelerator to improve memory utilisation and isolate the data of different tasks when local memory is sufficient, thereby reducing the time required for context switches.

### A. The Task Model

In our system, which consists of a CPU and a DNN accelerator, we consider a set of  $n$  sporadic tasks, denoted by  $\Gamma = \{\tau_1, \dots, \tau_n\}$ , and a shared resource (e.g., Gemmini<sup>RT</sup>) which may be required by a task. On startup, each task is initialised, followed by the requisite access of the shared resource until its execution completes. Each task  $\tau_i$  is characterised by a tuple of parameters  $(P_i, T_i, D_i, C_i^{LO}, C_i^{HI}, L_i, \eta_i)$ .  $P_i$  denotes the priority of task  $\tau_i$ , which is determined using a fixed priority assignment scheme in the system.  $T_i$  denotes the period or minimum inter-arrival time of task  $\tau_i$ .  $D_i$  denotes the deadline of task  $\tau_i$ .  $C_i^{LO}$  denotes the LO-WCET of task  $\tau_i$ , and  $C_i^{HI}$  denotes its HI-WCET, with the condition that for any task  $C_i^{LO} \leq C_i^{HI}$ .  $L_i$  denotes the criticality level of task  $\tau_i$ .  $\eta_i$  denotes the number of scratchpad banks allocated to task  $\tau_i$ .

### B. Worst-case Response Time Analysis

We now present the WCRT analysis. We begin by categorising tasks for a given task  $\tau_i$ . This categorisation is based on priority and criticality levels, leading to the formation of four categories:  $hpH(\tau_i)$  and  $hpL(\tau_i)$  denote the sets of HI-tasks and LO-tasks with higher priority than  $\tau_i$ ;  $lpH(\tau_i)$  and  $lpL(\tau_i)$  denote the set of HI-tasks and LO-tasks with lower priority than  $\tau_i$ .

Blocking occurs if a task is not scheduled when it is supposed to be, classified as either *criticality-inversion blocking* (*ci-blocking*), caused by a LO-task (even a high-priority LO-task) holding a shared resource needed by a HI-task in mode transition or HI-mode; or *priority-inversion blocking* (*pi-blocking*), caused by a task holding a resource needed by a higher-priority task.

In MESC, we only need to consider the blocking caused by tasks requiring the accelerator, as the context switches for CPU-only tasks are completed in such a short time that they cause



negligible blocking. We can initially perform a preliminary task partitioning using a function  $F(\Gamma)$ , which returns the subset of tasks within  $\Gamma$  that require accelerator utilisation. The complement of this subset, representing the tasks that only require the CPU, can be returned using  $\bar{F}(\Gamma)$ . For the accelerator, we have reduced the blocking impact caused by shared resources to the instruction level. This can be achieved by utilising a function  $I(\Gamma)$ , which provides the longest execution time of instructions on the accelerator within the task set  $\Gamma$ . The task scheduler operates at intervals of  $T_{sr}$  to routinely maintain system status and assess the necessity for context switches. We define some notations to describe the time consumption in different scenarios:  $\Upsilon_{Asr}^S$  and  $\Upsilon_{Asr}^R$  represent the maximum durations of the combined CPU context switch and accelerator context-saving time when a context save is required, and context-restoring time when a context restore is required, respectively.  $\Upsilon_{Csr}$  denotes the maximum duration of the CPU check time.  $\Upsilon_{Csr}^C$ <sup>8</sup> denotes the maximum time required for a context switch involving the CPU-only tasks. All parameters can be determined through experimental measurements. In our theoretical analysis, the periodic maintenance tasks performed by the task scheduler under these conditions are considered as the tasks with the highest priority. Given that these scenarios represent different branches of the same task scheduler behaviour, they will not coexist in the kernel. For clarity, we list these scenarios separately and do not categorise them under  $hpH(\tau_i)$ .

We consider three cases for the schedulability analysis [69]:

- 1) The schedulability of LO-mode by computing the WCRT of each task that is released and completed in LO-mode. In LO-mode, any task may experience *pi-blocking* from accelerator-required tasks in  $lpH(\tau_i)$  and  $lpL(\tau_i)$ . We use  $PB_i^{LO}$  and  $B_i^{LO}$  to respectively represent the maximum *pi-blocking* time and total blocking time task  $\tau_i$  may experience.
- 2) The schedulability of the HI-mode by computing the WCRT of each HI-task that is released and completed in HI-mode. In HI-mode, a HI-task may encounter *pi-blocking* due to accelerator-required tasks within both  $lpH(\tau_i)$  and  $lpL(\tau_i)$ , as well as *ci-blocking* resulting from accelerator-required tasks across  $lpL(\tau_i)$  and  $hpL(\tau_i)$ . We use  $PB_i^{HI}$ ,  $CB_i^{HI}$  and  $B_i^{HI}$  to respectively represent the maximum *pi-blocking* time, *ci-blocking* time and total blocking time that HI-task  $\tau_i$  may experience.
- 3) The schedulability of the mode transition by computing the WCRT of each HI-task that is released in either LO-mode or mode transition and completed in either HI-mode or mode transition. During mode transition, the blocking scenarios encountered are similar to those in HI-mode. We use  $PB_i^*$ ,  $CB_i^*$  and  $B_i^*$  to respectively represent the maximum *pi-blocking* time, *ci-blocking* time and total blocking time that HI-task  $\tau_i$  may experience.

<sup>8</sup>The context switch for CPU-only tasks has not been further differentiated into save and restore process because, unlike the accelerator, this duration is significantly shorter, and the time for saving and restoring is similar, lacking the distinct variability observed with accelerator.

**Response time analysis for LO-mode.** We consider a task  $\tau_i$  released and finished in LO-mode. The task  $\tau_i$  may experience at most one instance of *pi-blocking*. In our scheduling model, which involves periodic checks by the scheduler, if a higher-priority task arrives but misses the inspection window, it could face additional blocking equivalent to  $T_{sr}$  on top of the longest execution time of instructions from accelerator-required tasks in  $lpH(\tau_i)$  and  $lpL(\tau_i)$ . This duration is represented by Eq. 1.

$$PB_i^{LO} = I\left(F(lpH(\tau_i) \cup lpL(\tau_i))\right) + T_{sr} \quad (1)$$

As there is no criticality inversion in LO-mode,  $B_i^{LO}$  is equal to  $PB_i^{LO}$ , as shown in Eq. 2:

$$B_i^{LO} = PB_i^{LO} \quad (2)$$

In LO-mode, the response time  $R_i^{LO}$  for task  $\tau_i$  comprises the sum of its blocking time, execution time, preemption overhead, and the time preempted by other tasks, including the scheduler. The preemption overhead<sup>9</sup> is  $\Upsilon_{Asr}^S + \Upsilon_{Asr}^R$ , with the latter accounting for the potential additional time required if a higher-priority task arrives while a preempted task is restoring its context. The time during which the task is preempted by other tasks includes not only the execution time of those tasks but also the overhead of context save and restore (each preemption incurs one context save and one context restore). The  $R_i^{LO}$  is presented in Eq. 3.

$$\begin{aligned} R_i^{LO} = & B_i^{LO} + C_i^{LO} + \Upsilon_{Asr}^S + \Upsilon_{Asr}^R + \left\lceil \frac{R_i^{LO}}{T_{sr}} \right\rceil \Upsilon_{Csr} \\ & + \sum_{\tau_j \in \bar{F}(hpH(\tau_i) \cup hpL(\tau_i))} \left\lceil \frac{R_i^{LO}}{T_j} \right\rceil (2\Upsilon_{Csr}^C + C_j^{LO}) \\ & + \sum_{\tau_k \in F(hpH(\tau_i) \cup hpL(\tau_i))} \left\lceil \frac{R_i^{LO}}{T_k} \right\rceil (\Upsilon_{Asr}^S + \Upsilon_{Asr}^R + C_k^{LO}) \quad (3) \end{aligned}$$

**Response time analysis for HI-mode.** We consider a HI-task released and finished in HI-mode. In HI-mode, for a HI-task, all blocking caused by LO-tasks is treated as *ci-blocking*. Therefore, the *pi-blocking* experienced by a HI-task should only include blocking caused by lower-priority HI-tasks (Eq. 4).

$$PB_i^{HI} = I\left(F(lpH(\tau_i))\right) + T_{sr} \quad (4)$$

For  $CB_i^{HI}$ , since in HI-mode, a LO-task, regardless of its priority, could only start executing when there are no active HI-tasks (as detailed in Sec. IV), it is necessary to also consider  $hpL(\tau_i)$  in the context of *ci-blocking* for the HI-task  $\tau_i$ , as shown in Eq. 5.

$$CB_i^{HI} = I\left(F(lpL(\tau_i) \cup hpL(\tau_i))\right) + T_{sr} \quad (5)$$

The  $B_i^{HI}$  is the maximum of  $CB_i^{HI}$  and  $PB_i^{HI}$ .

$$B_i^{HI} = I\left(F(lpL(\tau_i) \cup hpL(\tau_i) \cup lpH(\tau_i))\right) + T_{sr} \quad (6)$$

<sup>9</sup>In fact, a more rigorous analysis should consider  $\max(\Upsilon_{Asr}^S + \Upsilon_{Asr}^R, 2\Upsilon_{Csr}^C)$ ; however, since context switching involving the accelerator necessarily includes the entire SoC's context switch. Thus in the worst case, the former strictly exceeds the latter.

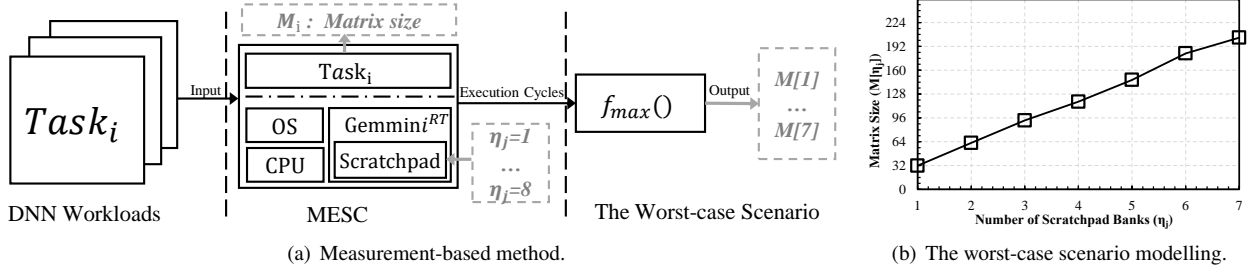


Fig. 6. Analysis of maximum input matrix sizes accommodated by different scratchpad capacities without impacting task execution speeds, based on varying workloads. The graph plots the sum of the sizes of the original input matrices against scratchpad capacities (the scratchpad consists of 8 banks, each is 32KB).

For HI-task  $\tau_i$ , the response time  $R_i^{\text{HI}}$ , as shown in Eq. 7, is composed of similar components in HI-mode as in LO-mode. Due to the scheduling preference for HI-tasks in HI-mode, only tasks from  $hpH(\tau_i)$  can preempt  $\tau_i$ .

$$\begin{aligned}
 R_i^{\text{HI}} &= B_i^{\text{HI}} + C_i^{\text{HI}} + \Upsilon_{\text{Asr}}^{\text{S}} + \Upsilon_{\text{Asr}}^{\text{R}} + \left\lceil \frac{R_i^{\text{HI}}}{T_{\text{Sr}}} \right\rceil \Upsilon_{\text{Csr}} \\
 &+ \sum_{\tau_j \in \bar{F}(hpH(\tau_i))} \left\lceil \frac{R_j^{\text{HI}}}{T_j} \right\rceil (2\Upsilon_{\text{Csr}}^{\text{C}} + C_j^{\text{HI}}) \\
 &+ \sum_{\tau_k \in F(hpH(\tau_i))} \left\lceil \frac{R_k^{\text{HI}}}{T_k} \right\rceil (\Upsilon_{\text{Asr}}^{\text{S}} + \Upsilon_{\text{Asr}}^{\text{R}} + C_k^{\text{HI}}) \quad (7)
 \end{aligned}$$

**Response time analysis for mode transition.** During the mode transition, for a HI-task, the experienced *pi-blocking* time, *ci-blocking* time and blocking time are equivalent to those in HI-mode, as shown in Eq. 8, 9 and 10.

$$PB_i^* = PB_i^{\text{HI}} \quad (8)$$

$$CB_i^* = CB_i^{\text{HI}} \quad (9)$$

$$B_i^* = B_i^{\text{HI}} \quad (10)$$

Given that any preemption of task  $\tau_i$  by LO-tasks can only occur in LO-mode, the response time  $R_i^*$  in mode transition is illustrated in Eq. 11.

$$\begin{aligned}
 R_i^* &= B_i^* + C_i^{\text{HI}} + \Upsilon_{\text{Asr}}^{\text{S}} + \Upsilon_{\text{Asr}}^{\text{R}} + \left\lceil \frac{R_i^*}{T_{\text{Sr}}} \right\rceil \Upsilon_{\text{Csr}} \\
 &+ \sum_{\tau_j \in \bar{F}(hpL(\tau_i))} \left\lceil \frac{R_j^{\text{LO}}}{T_j} \right\rceil (2\Upsilon_{\text{Csr}}^{\text{C}} + C_j^{\text{LO}}) \\
 &+ \sum_{\tau_k \in \bar{F}(hpH(\tau_i))} \left\lceil \frac{R_k^*}{T_k} \right\rceil (2\Upsilon_{\text{Csr}}^{\text{C}} + C_k^{\text{HI}}) \\
 &+ \sum_{\tau_m \in F(hpL(\tau_i))} \left\lceil \frac{R_m^{\text{LO}}}{T_m} \right\rceil (\Upsilon_{\text{Asr}}^{\text{S}} + \Upsilon_{\text{Asr}}^{\text{R}} + C_m^{\text{LO}}) \\
 &+ \sum_{\tau_n \in F(hpH(\tau_i))} \left\lceil \frac{R_n^*}{T_n} \right\rceil (\Upsilon_{\text{Asr}}^{\text{S}} + \Upsilon_{\text{Asr}}^{\text{R}} + C_n^{\text{HI}}) \quad (11)
 \end{aligned}$$

### C. Accelerator Local Memory Allocation Method

To facilitate the work of the address remapper (Sec. V), we need an analytical approach to illustrate the relationship between the size of task input matrices and accelerator local memory

utilisation. As pixel density or sequence data increases, the input matrices (or tensors) for DNNs expand, thereby elevating the demand for accelerator local memory and computational requirements. This offline analysis helps determine the minimum memory required to maintain task execution speeds without compromising performance.

We illustrate the mathematical relationship between the total size of task original input matrices and the minimal memory required ( $\eta_i$  for  $\tau_i$ ) under the constraint of unchanged execution times, as depicted in Fig. 6(a). Tests were conducted on all tasks mentioned in Sec. III, measuring their original input matrix sizes and execution times across different memory capacities. We then identified the minimal memory size that allows for optimal execution times (i.e., execution times when tasks are not constrained in their use of system resources) without exceeding a set execution threshold, as shown in Fig. 6(b). In matrix computations, the generation of a substantial volume of intermediate results typically elevates the demand for accelerator local memory. In Gemmini<sup>RT</sup>, however, these intermediate results are stored separately in accumulator, distinct from the primary data storage, which helps mitigate their impact on memory requirements. Our modelling is focused on scratchpad. Thus, the correlation between the size of input matrices and the required memory size is remarkably close.

Despite this specific arrangement in Gemmini<sup>RT</sup>, our objective is to establish a universal methodology, not solely applicable to this particular architecture. For accelerators where intermediate results are intermixed with the original data, we can apply a similar modelling approach. By considering the worst-case scenario, we model both data types together to determine the optimal local memory allocation within the accelerator to ensure efficient processing across different scenarios.

## VIII. EXPERIMENTAL EVALUATION

**Experimental platform.** We built the MESC on an AMD Alveo U280 evaluation board, utilising a Rocket core [68] as the RISC-V processor. The processor was instantiated with a 5-stage pipeline and single-width dispatch, featuring a 32 KB L1 I-cache and a 32 KB L1 D-cache, along with a shared 512 KB L2 cache, 4 GB external memory. In the Gemmini<sup>RT</sup> component, we adhered closely to the default configuration provided by the open-source Gemmini project. This configuration includes a 256 KB scratchpad, a 64 KB accumulator, and a systolic array

with a single tile consisting of 256 PEs. The system supports a maximum transfer size of 64 bytes with a bus width of 128 for the DMA. The only modification we introduced is an increase in the number of scratchpad banks from 4 to 8. The software stack, including the OS kernel and workloads, was compiled using the RISC-V GNU toolchain. We adopted FreeRTOS (v.10.5) as the OS, incorporating the modifications detailed in Sec. VI. The system runs at 100 MHz.

**Task set setup.** To ensure a comprehensive evaluation of MESC, we utilise the workloads from the experiments described in Sec. III — the DNN workloads and test files developed by UC Berkeley for Gemmini [15], [16]. Tasks are randomly selected based on a uniform distribution to form each task set. The task set parameters used were:

- Task utilisations  $U_i$  were generated using UUnifast [70], providing an unbiased distribution.
- The  $C_i^{LO}$  for each task was directly obtained by retrieving relevant task information from the dataset. The  $C_i^{HI}$  was then calculated by  $C_i^{HI} = CF \cdot C_i^{LO}$ , where  $CF$  was the criticality factor, default  $CF = 2.0$  [66], [71].
- The task periods  $T_i$  were determined by  $T_i = \frac{C_i^{LO}}{U_i}$ .
- Tasks were assigned implicit deadlines, implying  $D_i = T_i$ .
- Tasks were managed using fixed-priority scheduling, with priorities assigned in descending order of  $T_i$ .
- The proportion of HI-tasks within the task set was determined by the criticality proportion  $\gamma$ , default  $\gamma = 0.5$ .

In our schedulability tests, the total utilisation of the task set was set from 0.5 to 0.95, with an incremental step of 0.1. For each utilisation level, 1000 task sets were generated. The number of tasks in each set was determined by  $\beta$ , default  $\beta = 10$ . The operational interval of the scheduler was  $T_{sr} = 5000$  cycles.

#### A. Context-switching Overhead and Blocking Duration

**Experimental setup.** We initialised the MESC in LO-mode and introduced random disturbances for HI-tasks by randomly increasing the input of some HI-tasks before execution, thereby extending their execution times beyond  $C_i^{LO}$  to prompt mode switch of the system. After running the MESC for a set duration, we recorded the time taken for context save, context restore, priority inversions, and criticality inversions. Furthermore, we conducted experiments with the bank allocation strategy and the context-switching mechanism disabled in the MESC, aiming to determine the acceleration ratios provided by each component.

**Obs. 1.** By comparing the complete MESC with the system lacking bank allocations, in terms of context save and restore durations in Fig. 7, we observe an increase in context-switching times by 4000 to 6000 cycles when the bank allocation is removed. This increased results from the average case, where the address remapper and bank allocation method reduced the frequency and size of data transfers within the accelerator, achieving an acceleration of approximately 20% to 30%.

**Obs. 2.** Comparing the duration of priority inversion and criticality inversion in the complete MESC in Fig. 7, we observed that, on average, the duration of criticality inversion was reduced by 3000 to 5000 cycles compared to priority

inversion. This reduction was due to the system rules in HI-mode (detailed in Sec. IV). When a HI-task is preempting a LO-task, only the data of the currently running LO-task remains in the scratchpad, significantly increasing the likelihood that the accelerator’s local memory is sufficient during the context switch, thus eliminating the need to handle computation data in the scratchpad and then accelerating the context-switching process. Further comparisons of these durations with a system where context switch functionality was removed revealed significant improvements. For priority inversion durations, we achieved an acceleration exceeding 250 times; for criticality inversion durations, the acceleration surpassed 300 times. This substantial increase in efficiency stems from the fact that without context switch, the accelerator, being a shared resource with long critical-section durations, has to wait until the previous task has completely finished, leading to a substantial waiting period. We effectively resolved this issue at the instruction level, allowing timely task transitions and significantly minimising waiting times.

#### B. Successful Ratio

**Experimental setup.** We initiated MESC in LO-mode and after an initial set running period, recorded the successful ratio (i.e., the proportion of runs with no tasks missing their deadlines during that period) of task sets in LO-mode, mode transition and HI-mode. To explore the impact of context switch on successful ratio, we conducted experiments on the MESC without the context-switching mechanism. We also introduced Adaptive Mixed Criticality (AMC) [27] with a strategy of terminating all LO-tasks in HI-mode for comparison.

**Obs. 3.** Fig. 8 shows the successful ratio of MESC under different modes in the default conditions. As observed, even at a utilisation rate of 0.95, MESC maintains considerable schedulability. However, when the context-switching mechanism of MESC is disabled, its success rate plummets to 0 at a utilisation rate of 0.85. This highlights the crucial impact of the context-switching mechanism on successful ratio. Additionally, we included comparisons with AMC systems, both with and without context switch. It is evident that without the context-switching mechanism, even terminating all LO-tasks in HI-mode does not ensure system success. Furthermore, without the context-switching mechanism, AMC shows significantly higher success rates than MESC. However, when both systems incorporate context switch, AMC only marginally outperforms MESC in scheduling success. This could be attributed to the fact that although MESC maintains the execution of LO-tasks in HI-mode, the context-switching mechanism enables greater flexibility and real-time responsiveness.

#### C. Successful Ratio in HI-mode under Various Variables

**Experimental setup.** Using the same experimental setup as the successful ratio test, we examined the impact of  $\gamma$  and  $\beta$  on the successful ratio of HI-tasks in the HI-mode.

**Obs. 4.** As shown in Fig. 9(a), when  $\gamma$  increases, the successful ratio of the system significantly decreases in a stable trend. This observation is caused by the higher proportion of HI-tasks, which

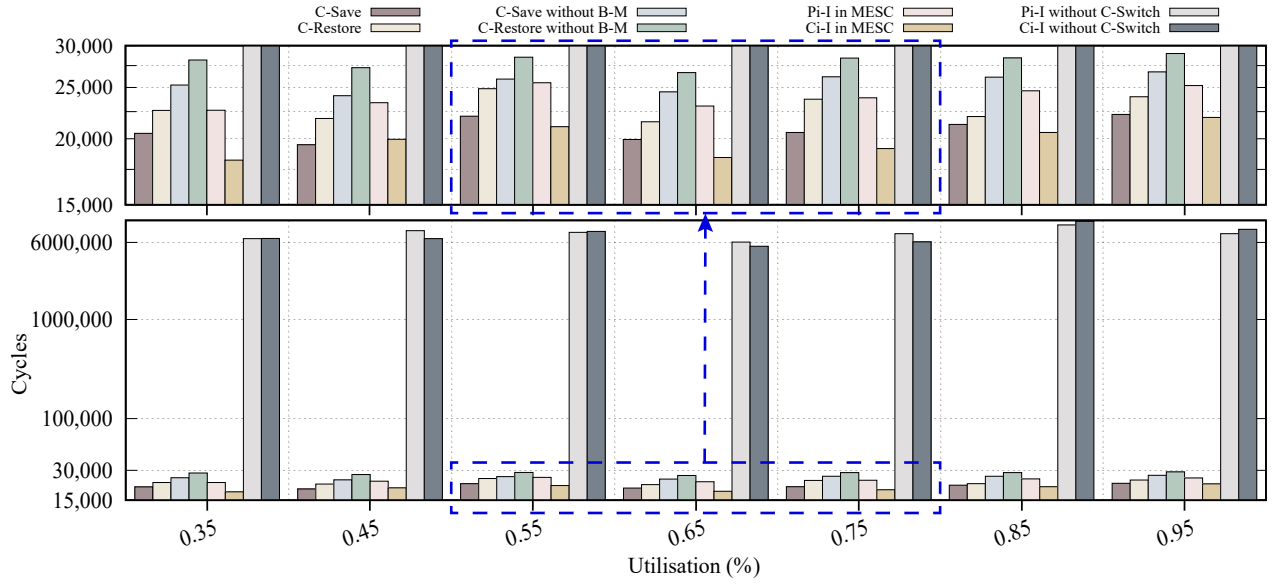


Fig. 7. Comparison of the impact of different components of MESC under various utilizations. “C-Save” and “C-Restore” represent the cycles for context saving and restoring. “C-Save without B-M” and “C-Restore without B-M” show the cycles without the bank model. “Pi-I in MESC” and “Ci-I in MESC” denote the cycles for priority inversion and criticality inversion within MESC, respectively. “Pi-I without C-Switch” and “Ci-I without C-Switch” indicate the cycles for priority inversion and criticality inversion without the context-switching mechanism.

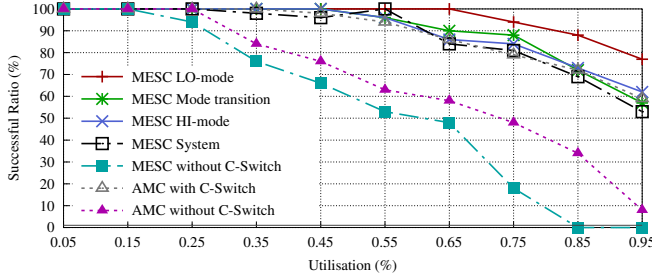


Fig. 8. Successful ratio of MESC in different modes, as well as the successful ratio of the entire system with and without context-switching functionality, compared to the successful ratio of the system with and without context switching when applying AMC. “C-Switch” denotes context switch.

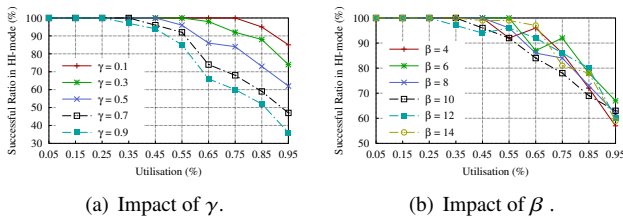


Fig. 9. Successful ratio of MESC in HI-mode under various conditions.

increases the system load in HI-mode and intensifies resource contention for the accelerator, thereby reducing the likelihood of all tasks meeting their deadlines. In contrast, as shown in Fig. 9(b), the decrease in the successful ratio is minimal as  $\beta$  decreases in MESC. This is because changes in the number of tasks do not fundamentally increase the system load, and thus do not impact system schedulability as directly as an increase in  $\gamma$ .

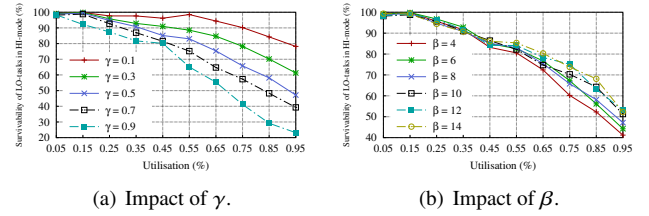


Fig. 10. Survivability of LO-tasks in HI-mode under various conditions.

Additionally, this stability is partially attributed to the effective context switching, which ensures that the system can promptly handle arriving high-priority tasks regardless of changes in  $\beta$ , allowing the system to maintain a stable success rate even as the number of tasks decreases.

#### D. Survivability of LO-tasks under Various Variables

We demonstrate the survivability [34] of LO-tasks in HI-mode by running the workloads mentioned in Sec. III using the MESC. Survivability is defined as the ratio of the number of LO-tasks that complete to the number of LO-tasks released in HI-mode.

**Experimental setup.** Using the same experimental setup as in the successful ratio test, we examined the impact of  $\gamma$  and  $\beta$  on the survivability of LO-tasks in HI-mode.

**Obs. 5.** Fig. 10 illustrates the impact of  $\gamma$  and  $\beta$  on the survivability of LO-tasks in HI-mode. As shown in Fig. 10(a), the survivability of LO-tasks significantly decreases with an increase in  $\gamma$ , as expected. However, even in the most extreme cases, we still maintain over 20% survivability of LO-tasks, which is crucial for the stability of the system. Similar to the successful ratio tests, as shown in Fig. 10(b), the impact of  $\beta$  on the

survivability of LO-tasks remains very low. This can be attributed to the context-switching mechanism that allows LO-tasks to run opportunistically, alleviating the pressure on their survivability across different  $\beta$ , especially at lower utilisation levels.

### E. Hardware Overhead

**Experimental setup.** We compared the hardware overhead and power consumption of Gemmini<sup>RT</sup> with the default version of Gemmini. We examined the overhead distribution across its various components, including the controller units, systolic array, reservation station, and scratchpad. The remapping block in the scratchpad is configured as 4KB. We have synthesised and implemented Gemmini<sup>RT</sup> on the Alveo U280 FPGA using Vivado (v2021.1) to examine the consumption. The hardware consumption is mainly summarised as the usage of Look-Up-Tables (LUTs), registers, DSPs, RAMs and power consumption.

**Obs. 6.** The resource efficiency of Gemmini<sup>RT</sup> is demonstrated in Tbl. II. Note that Gemmini<sup>RT</sup> only introduced an additional 8815 (4.8%) LUTs, 2861 (4.3%) registers, 1 (0.5%) DSPs, 4KB (1.2%) RAMs, and 61 (5.2%) mW. Although this comparison is based on Gemmini’s open-source default configuration, it is important to note that this configuration includes only a very small systolic array (with a single tile), only 256KB of scratchpad and 64KB of accumulator. In practical scenarios, increasing these configurations would not add to the overhead of the context-switching functionality we introduced, and the relative overhead would likely be even lower. This is because, even with only a single tile, the overhead and power consumption of the systolic array (64.8% LUTs, 42.1% registers, and 45.5% mW) still constitutes a major portion of the overall hardware costs.

## IX. RELATED WORK

The problems of priority and/or criticality inversions in co-processors have received considerable attention in the research community. Some researchers adopt scheduling protocols to mitigate this problem. For example, GPU<sub>sync</sub> [21] uses predictive scheduling to allow task status to be transferred among different GPUs. Similarly, Wu et al. [72] optimise resource utilisation by managing the execution of computation graph subgraphs and migrating tasks across different system components. Whilst this scheduling strategy addresses these problems to some extent, it overly relies on hardware resources and struggles to manage scenarios with high task density. Park et al. [73] propose Chimaera, which abandons the current executing task to free up resources for higher-priority tasks. However, this approach imposes overly strict constraints for real-time application scenarios. Some researchers have attempted to address these issues at the software level. For instance, Liu et al. [17]–[20] partition the software into smaller chunks. Similarly, [32], [49], [74], [75] adopt approaches that break tasks into thread-level segments. Although these methods offer improvements, they typically involve limited preemption and necessitate substantial modifications to the software workloads. Hartmann and Margull [76] introduce a software-based preemption technique that enables controllable

TABLE II  
HARDWARE OVERHEAD (IMPLEMENTED ON FPGA)

	LUTS	Registers	DSP	BRAM (KB)	Power (mW)
<b>Gemmini</b>	182535	65317	207	320	1171
<b>Gemmini<sup>RT</sup></b>	191350	68178	208	324	1232
<b>Additional Overhead</b>	8815 (4.8%)	2861 (4.3%)	1 (0.5%)	4 (1.2%)	61 (5.2%)

TABLE III  
DECOMPOSED OVERHEAD OF EACH COMPONENT IN THE GEMMINI<sup>RT</sup>

Gemmini <sup>RT</sup>	LUTS (%)	Registers (%)	DSP (%)	BRAM (%)	Power (%)
<b>Controller</b>	4.8%	5.3%	3.4%	0	3.1%
<b>Systolic Array</b>	64.8%	42.1%	0	0	45.5%
<b>Reservation Station</b>	4.1%	10.1%	2.4%	0	2.8%
<b>Scratchpad</b>	20.2%	32.4%	20.2%	100%	36.6%
<b>Other Module</b>	6.1%	10.1%	74.0%	0	12.0%

preemption at the instruction level, but there is still significant room for improvement in terms of general applicability. Lee et al. [77] implemented fine-grained context switching on GPUs by combining workload partitioning with scheduling, which significantly alleviates the problem. However, this approach incurs a high cost, as it requires rolling back and re-executing previously preempted subtasks. Some researchers address these issues at the hardware level. Jiang et al. [64] propose the MCS-I/OV to support context switch for co-processors, but this solution is limited to I/O. Tanasic et al. [78] and Sasongko et al. [79] create hardware enhancements that facilitate preemption by enabling context switch at the thread level. Moreover, most industrial embedded GPUs, such as those in the Jetson TX2 (Pascal) [80]–[82], Xavier (Volta) [83], and Orin (Ampere) [84] with compute capability below 6.0, also lack the capability for a fine-grained context switch. Although starting with the Pascal architecture [85] (compute capability 6.0 and higher), these GPUs support instruction-level preemption within compute-type tasks, this granularity still does not fully meet the needs of all tasks requiring a fine-grained context switch.

## X. CONCLUSION

In this paper we present MESC, a systematic framework that incorporates hardware and software solutions, forming a comprehensive solution for heterogeneous MCSs. We illustrate how instruction-level context switch for DNN accelerators can be enabled under MESC. Through theoretical analysis and experiments, we show that compared to traditional non-preemptive accelerators, the MESC framework achieves an improvement of 2 orders of magnitude in resolving algorithmic priority and criticality inversions. Additionally, MESC significantly enhances the system’s timing performance and its ability to manage high-priority and HI-tasks. It also improves the survivability of LO-tasks in HI-mode, whilst incurring minimal hardware overhead.

## XI. ACKNOWLEDGEMENT

The authors would like to thank the anonymous reviewers and the shepherd for their invaluable feedback on this paper. This work was supported by the National Natural Science Foundation of China (Grant No. 62472086), the Natural Science Foundation of Jiangsu Province (Grants No. BK20243042), and the Start-up Research Fund of Southeast University (Grant No. RF1028624005). In his first work, Jiapeng Guan would like to express his deepest gratitude to his family, with heartfelt wishes for his father's swift recovery and a healthy life free from pain in the days ahead. Additionally, Jiapeng extends his thanks to Yuanfeng Xu for the supports during system development.

## REFERENCES

- [1] A. Burns, R. Davis, Mixed criticality systems—a review, Department of Computer Science, University of York, Tech. Rep (2013) 1–69.
- [2] Z. Jiang, S. Zhao, R. Wei, D. Yang, R. Paterson, N. Guan, Y. Zhuang, N. C. Audsley, Bridging the pragmatic gaps for mixed-criticality systems in the automotive industry, *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems (TCAD)* 41 (4) (2021) 1116–1129.
- [3] A. Naghavi, S. Safari, S. Hessabi, Tolerating permanent faults with low-energy overhead in multicore mixed-criticality systems, *IEEE Transactions on Emerging Topics in Computing (TETC)* 10 (2) (2021) 985–996.
- [4] C. Hernández, J. Flich, R. Paredes, C.-A. Lefebvre, I. Allende, J. Abella, D. Trilla, M. Matschnig, B. Fischer, K. Schwarz, et al., Selene: Self-monitored dependable platform for high-performance safety-critical systems, in: 2020 23rd Euromicro Conference on Digital System Design (DSD), IEEE, 2020, pp. 370–377.
- [5] Z. Hu, J. Luo, X. Fang, K. Xiao, B. Hu, L. Chen, Real-time schedule algorithm with temporal and spatial isolation feature for mixed criticality system, in: 2021 7th International Symposium on System and Software Reliability (ISSSR), IEEE, 2021, pp. 99–108.
- [6] D. Yang, X. Li, X. Dai, R. Zhang, L. Qi, W. Zhang, Z. Jiang, All in one network for driver attention monitoring, in: ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE, 2020, pp. 2258–2262.
- [7] Z. Jiang, S. Zhao, P. Dong, D. Yang, R. Wei, N. Guan, N. Audsley, Rethinking mixed-criticality architecture for automotive industry, in: 2020 IEEE 38th International Conference on Computer Design (ICCD), IEEE, 2020, pp. 510–517.
- [8] I. ISO, 26262: Road vehicles—functional safety (2018).
- [9] Z. Jiang, X. Dai, P. Dong, R. Wei, D. Yang, N. C. Audsley, N. Guan, Toward an analysable, scalable, energy-efficient i/o virtualization for mixed-criticality systems, *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems (TCAD)* (2022).
- [10] D. Ottaviano, M. Cinque, G. Manduchi, S. Dubbioso, Virtualization of accelerators in embedded systems for mixed-criticality: Rpu exploitation for fusion diagnostics and control, *Fusion Engineering and Design (FED)* 190 (2023) 113518.
- [11] S. Majumder, J. F. D. Nielsen, A. la Cour-Harbo, H. Schiøler, T. Bak, A real-time on-chip network architecture for mixed criticality aerospace systems, *The Aeronautical Journal (Aeronaut. J)* 123 (1269) (2019) 1788–1806.
- [12] C. Latotzke, T. Gemmeke, Efficiency versus accuracy: A review of design techniques for dnn hardware accelerators, *IEEE Access* 9 (2021) 9785–9799.
- [13] S. Lu, M. Wang, S. Liang, J. Lin, Z. Wang, Hardware accelerator for multi-head attention and position-wise feed-forward in the transformer, in: 2020 IEEE 33rd International System-on-Chip Conference (SOCC), IEEE, 2020, pp. 84–89.
- [14] C. Fang, S. Guo, W. Wu, J. Lin, Z. Wang, M. K. Hsu, L. Liu, An efficient hardware accelerator for sparse transformer neural networks, in: 2022 IEEE International Symposium on Circuits and Systems (ISCAS), IEEE, 2022, pp. 2670–2674.
- [15] H. Genc, A. Haj-Ali, V. Iyer, A. Amid, H. Mao, J. Wright, C. Schmidt, J. Zhao, A. Ou, M. Banister, et al., Gemmini: An agile systolic array generator enabling systematic evaluations of deep-learning architectures, *arXiv preprint arXiv: 1911.09925* 3 (2019) 25.
- [16] H. Genc, S. Kim, A. Amid, A. Haj-Ali, V. Iyer, P. Prakash, J. Zhao, D. Grubb, H. Liew, H. Mao, et al., Gemmini: Enabling systematic deep-learning architecture evaluation via full-stack integration, in: 2021 58th ACM/IEEE Design Automation Conference (DAC), IEEE, 2021, pp. 769–774.
- [17] S. Liu, S. Yao, X. Fu, R. Tabish, S. Yu, A. Bansal, H. Yun, L. Sha, T. Abdelzaher, On removing algorithmic priority inversion from mission-critical machine inference pipelines, in: 2020 IEEE Real-Time Systems Symposium (RTSS), IEEE, 2020, pp. 319–332.
- [18] S. Liu, L. Sha, T. Abdelzaher, Criticality-based data segmentation and resource allocation in machine, *Artificial Intelligence for Edge Computing (AI4EC)* (2023) 335.
- [19] S. Liu, L. Sha, T. Abdelzaher, Criticality-based data segmentation and resource allocation in machine inference pipelines, in: *Artificial Intelligence for Edge Computing (AI4EC)*, Springer, 2023, pp. 335–352.
- [20] S. Liu, S. Yao, X. Fu, R. Tabish, S. Yu, A. Bansal, H. Yun, L. Sha, T. Abdelzaher, Taming algorithmic priority inversion in mission-critical perception pipelines, *Communications of the ACM (Commun. ACM)* 67 (2) (2024) 110–117.
- [21] G. A. Elliott, B. C. Ward, J. H. Anderson, Gpufreq: A framework for real-time gpu management, in: 2013 IEEE 34th Real-Time Systems Symposium (RTSS), IEEE, 2013, pp. 33–44.
- [22] R. Ernst, M. Di Natale, Mixed criticality systems—a history of misconceptions?, *IEEE Design & Test* 33 (5) (2016) 65–74.
- [23] P. Huang, P. Kumar, N. Stoimenov, L. Thiele, Interference constraint graph—a new specification for mixed-criticality systems, in: 2013 IEEE 18th Conference on Emerging Technologies & Factory Automation (ETF), IEEE, 2013, pp. 1–8.
- [24] M. Paulitsch, O. M. Duarte, H. Karray, K. Mueller, D. Muench, J. Nowotch, Mixed-criticality embedded systems—a balance ensuring partitioning and performance, in: 2015 Euromicro Conference on Digital System Design (DSD), IEEE, 2015, pp. 453–461.
- [25] R. I. Davis, A. Burns, I. Bate, Compensating adaptive mixed criticality scheduling, in: *Proceedings of the 30th International Conference on Real-Time Networks and Systems (RTNS)*, 2022, pp. 81–93.
- [26] L. Zeng, C. Xu, R. Li, Partition and scheduling of the mixed-criticality tasks based on probability, *IEEE Access* 7 (2019) 87837–87848.
- [27] S. K. Baruah, A. Burns, R. I. Davis, Response-time analysis for mixed criticality systems, in: 2011 IEEE 32nd Real-Time Systems Symposium (RTSS), IEEE, 2011, pp. 34–43.
- [28] A. Burns, S. Baruah, Towards a more practical model for mixed criticality systems, in: *Workshop on Mixed-Criticality Systems (WMC)*, 2013.
- [29] J. Lee, M. Kim, Generalized models of mixed-criticality systems for real-time scheduling, *IEEE Transactions on Engineering and Computer Science (Trans. Eng. Comput. Sci.)* 1 (1-50) (2020) 51.
- [30] M.-K. Yoon, J.-E. Kim, R. Bradford, Z. Shao, Timedice: Schedulability-preserving priority inversion for mitigating covert timing channels between real-time partitions, in: 2022 52nd Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN), IEEE, 2022, pp. 453–465.
- [31] N. C. Gaitan, I. Zagan, V. G. Gaitan, Predictable cpu architecture designed for small real-time application—concept and theory of operation, *International Journal of Advanced Computer Science and Applications (IJACSA)* 6 (4) (2015).
- [32] H. Zhou, G. Tong, C. Liu, Gpes: A preemptive execution system for gpgpu computing, in: 21st IEEE Real-Time and Embedded Technology and Applications Symposium (RTAS), IEEE, 2015, pp. 87–97.
- [33] K. Suo, Y. Shi, C.-C. Hung, P. Bobbie, Quantifying context switch overhead of artificial intelligence workloads on the cloud and edges, in: *Proceedings of the 36th Annual ACM Symposium on Applied Computing (SAC)*, 2021, pp. 1182–1189.
- [34] Z. Jiang, X. Dai, N. Audsley, Hiart-mcs: High resilience and approximated computing architecture for imprecise mixed-criticality systems, in: 2021 IEEE Real-Time Systems Symposium (RTSS), IEEE, 2021, pp. 290–303.
- [35] S. Baruah, A. Burns, Fixed-priority scheduling of dual-criticality systems, in: *Proceedings of the 21st International Conference on Real-Time Networks and Systems (RTNS)*, 2013, pp. 173–181.
- [36] Z. Li, S. He, Fixed-priority scheduling for two-phase mixed-criticality systems, *ACM Transactions on Embedded Computing Systems (TECS)* 17 (2) (2017) 1–20.
- [37] A. Burns, R. I. Davis, Response time analysis for mixed criticality systems with arbitrary deadlines, in: 5th International Workshop on Mixed-Criticality Systems (WMC), York, 2017.



- [38] A. Burns, R. I. Davis, A survey of research into mixed criticality systems, *ACM Computing Surveys (CSUR)* 50 (6) (2017) 1–37.
- [39] D. Succi, P. Poplavko, S. Bensalem, M. Bozga, Mixed critical earliest deadline first, in: 2013 25th Euromicro Conference on Real-Time Systems (ECRTS), IEEE, 2013, pp. 93–102.
- [40] Q. Zhao, M. Qu, Z. Gu, H. Zeng, Minimizing stack memory for partitioned mixed-criticality scheduling on multiprocessor platforms, *ACM Transactions on Embedded Computing Systems (TECS)* 21 (2) (2022) 1–30.
- [41] Y.-W. Zhang, R.-K. Chen, Z. Gu, Energy-aware partitioned scheduling of imprecise mixed-criticality systems, *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems (TCAD)* (2023).
- [42] N. Chen, S. Zhao, I. Gray, A. Burns, S. Ji, W. Chang, Msrp-ft: Reliable resource sharing on multiprocessor mixed-criticality systems, in: 2022 IEEE 28th Real-Time and Embedded Technology and Applications Symposium (RTAS), IEEE, 2022, pp. 201–213.
- [43] Z. Dong, C. Liu, Schedulability analysis for coscheduling real-time tasks on multiprocessors, *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems (TCAD)* 41 (11) (2022) 4721–4732.
- [44] Y. Xiang, H. Kim, Pipelined data-parallel cpu/gpu scheduling for multi-dnn real-time inference, in: 2019 IEEE Real-Time Systems Symposium (RTSS), IEEE, 2019, pp. 392–405.
- [45] X. Lin, R. Liu, J. Xie, Q. Wei, Z. Zhou, X. Chen, Z. Huang, G. Lu, Online scheduling of cpu-npu co-inference for edge ai tasks, in: 2023 IEEE Wireless Communications and Networking Conference (WCNC), IEEE, 2023, pp. 1–6.
- [46] Z. Xu, D. Yang, C. Yin, J. Tang, Y. Wang, G. Xue, A co-scheduling framework for dnn models on mobile and edge devices with heterogeneous hardware, *IEEE Transactions on Mobile Computing (TMC)* 22 (3) (2021) 1275–1288.
- [47] Y. Wang, M. Karimi, Y. Xiang, H. Kim, Balancing energy efficiency and real-time performance in gpu scheduling, in: 2021 IEEE Real-Time Systems Symposium (RTSS), IEEE, 2021, pp. 110–122.
- [48] T. Amert, C. Nemitz, Work-in-progress: Impacts of critical-section granularity when accessing shared resources, in: 2023 IEEE Real-Time Systems Symposium (RTSS), IEEE, 2023, pp. 439–442.
- [49] C. Basaran, K.-D. Kang, Supporting preemptive task executions and memory copies in gpgpus, in: 2012 24th Euromicro Conference on Real-Time Systems (ECRTS), IEEE, 2012, pp. 287–296.
- [50] S. Vestal, Preemptive scheduling of multi-criticality systems with varying degrees of execution time assurance, in: 28th IEEE International Real-Time Systems Symposium (RTSS), IEEE, 2007, pp. 239–243.
- [51] S. Baruah, V. Bonifaci, G. D'Angelo, H. Li, A. Marchetti-Spaccamela, S. Van Der Ster, L. Stougie, The preemptive uniprocessor scheduling of mixed-criticality implicit-deadline sporadic task systems, in: 2012 24th Euromicro Conference on Real-Time Systems (ECRTS), IEEE, 2012, pp. 145–154.
- [52] Z. Guo, K. Yang, S. Vaidhun, S. Arefin, S. K. Das, H. Xiong, Uniprocessor mixed-criticality scheduling with graceful degradation by completion rate, in: 2018 IEEE Real-Time Systems Symposium (RTSS), IEEE, 2018, pp. 373–383.
- [53] Z. Jiang, K. Yang, N. Fisher, N. Audsley, Z. Dong, Pythia-mcs: Enabling quarter-clairvoyance in i/o-driven mixed-criticality systems, in: 2020 IEEE Real-Time Systems Symposium (RTSS), IEEE, 2020, pp. 38–50.
- [54] Z. Jiang, X. Dai, A. Burns, N. Audsley, Z. Gu, I. Gray, A high-resilience imprecise computing architecture for mixed-criticality systems, *IEEE Transactions on Computers (TC)* 72 (1) (2022) 29–42.
- [55] A. Amid, D. Biancolin, A. Gonzalez, D. Grubb, S. Karandikar, H. Liew, A. Magyar, H. Mao, A. Ou, N. Pemberton, et al., Chipyard—an integrated soc research and implementation environment, in: 2020 57th ACM/IEEE Design Automation Conference (DAC), IEEE, 2020, pp. 1–6.
- [56] M. Gao, J. Pu, X. Yang, M. Horowitz, C. Kozyrakakis, Tetris: Scalable and efficient neural network acceleration with 3d memory, in: Proceedings of the Twenty-Second International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS), 2017, pp. 751–764.
- [57] E. Qin, A. Samajdar, H. Kwon, V. Nadella, S. Srinivasan, D. Das, B. Kaul, T. Krishna, Sigma: A sparse and irregular gemm accelerator with flexible interconnects for dnn training, in: 2020 IEEE International Symposium on High Performance Computer Architecture (HPCA), IEEE, 2020, pp. 58–70.
- [58] D. Shin, J. Lee, J. Lee, H.-J. Yoo, Dnpu: An energy-efficient deep-learning processor with heterogeneous multi-core architecture, *IEEE Micro* 38 (5) (2018) 85–93.
- [59] D. A. N. Gookyi, E. Lee, K. Kim, S.-J. Jang, S.-S. Lee, Deep learning accelerators' configuration space exploration effect on performance and resource utilization: A gemmini case study, *Sensors* 23 (5) (2023) 2380.
- [60] F. N. Peccia, O. Bringmann, Integration of a systolic array based hardware accelerator into a dnn operator auto-tuning framework, *arXiv preprint arXiv: 2212.03034* (2022).
- [61] M. A. Shafique, K. Inayat, J. A. Lee, Csa based radix-4 gemmini systolic array for machine learning applications, in: Proceedings of the 13th International Symposium on Highly Efficient Accelerators and Reconfigurable Technologies (HEART), 2023, pp. 93–99.
- [62] J. Vieira, N. Roma, G. Falcao, P. Tomás, Gem5-accel: A pre-rtl simulation toolchain for accelerator architecture validation, *IEEE Computer Architecture Letters (CAL)* (2023).
- [63] J. Zou, X. Dai, J. A. McDermid, Context-aware graceful degradation for mixed-criticality scheduling in autonomous systems, *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems (TCAD)* (2023).
- [64] Z. Jiang, N. Audsley, P. Dong, N. Guan, X. Dai, L. Wei, Mcs-iov: Real-time i/o virtualization for mixed-criticality systems, in: 2019 IEEE Real-Time Systems Symposium (RTSS), IEEE, 2019, pp. 326–338.
- [65] Z. Jiang, P. Dong, R. Wei, Q. Zhao, Y. Wang, D. Zhu, Y. Zhuang, N. Audsley, Pspys: A time-predictable mixed-criticality system architecture based on arm trustzone, *Journal of Systems Architecture (JSA)* 123 (2022) 102368.
- [66] A. Burns, R. I. Davis, Schedulability analysis for adaptive mixed criticality systems with arbitrary deadlines and semi-clairvoyance, in: 2020 IEEE Real-Time Systems Symposium (RTSS), IEEE, 2020, pp. 12–24.
- [67] A. Burns, C. Jones, An approach to formally specifying the behaviour of mixed-criticality systems, in: Euromicro Conference on Real-Time Systems (ECRTS), ACM, 2022, p. 12.
- [68] K. Asanovic, R. Avizienis, J. Bachrach, S. Beamer, D. Biancolin, C. Celio, H. Cook, D. Dabbelt, J. Hauser, A. Izraelvitz, et al., The rocket chip generator, EECs Department, University of California, Berkeley, Tech. Rep. UCB/EECS-2016-17 4 (2016) 6–2.
- [69] S. K. Baruah, A. Burns, R. I. Davis, Response-time analysis for mixed criticality systems, in: IEEE International Real-Time Systems Symposium (RTSS), 2011.
- [70] E. Bini, G. C. Buttazzo, Measuring the performance of schedulability tests, *Real-Time Systems (Real-Time Syst.)* 30 (1–2) (2005) 129–154.
- [71] I. Bate, A. Burns, R. I. Davis, Analysis-runtime co-design for adaptive mixed criticality scheduling, in: 2022 IEEE 28th Real-Time and Embedded Technology and Applications Symposium (RTAS), IEEE, 2022, pp. 187–200.
- [72] X. Wu, J. Rao, W. Chen, H. Huang, C. Ding, H. Huang, Switchflow: preemptive multitasking for deep learning, in: Proceedings of the 22nd International Middleware Conference (IMC), 2021, pp. 146–158.
- [73] J. J. K. Park, Y. Park, S. Mahlke, Chimera: Collaborative preemption for multitasking on a shared gpu, *ACM SIGARCH Computer Architecture News (ACM SIGARCH Comput. Archit. News)* 43 (1) (2015) 593–606.
- [74] B. Wu, X. Liu, X. Zhou, C. Jiang, Flep: Enabling flexible and efficient preemption on gpus, *ACM SIGPLAN Notices* 52 (4) (2017) 483–496.
- [75] T. Amert, Z. Tong, S. Voronov, J. Bakita, F. D. Smith, J. H. Anderson, Timewall: Enabling time partitioning for real-time multicore+ accelerator platforms, in: 2021 IEEE Real-Time Systems Symposium (RTSS), IEEE, 2021, pp. 455–468.
- [76] C. Hartmann, U. Margull, Gpuart—an application-based limited preemptive gpu real-time scheduler for embedded systems, *Journal of Systems Architecture (JSA)* 97 (2019) 304–319.
- [77] H. Lee, H. Kim, C. Kim, H. Han, E. Seo, Idempotence-based preemptive gpu kernel scheduling for embedded systems, *IEEE Transactions on Computers (TC)* 70 (3) (2020) 332–346.
- [78] I. Tanasic, I. Gelado, J. Cabezas, A. Ramirez, N. Navarro, M. Valero, Enabling preemptive multiprogramming on gpus, in: Proceeding of the 41st Annual International Symposium on Computer Architecture (ISCA), IEEE, Minneapolis, Minnesota, USA, 2014, pp. 193–204.
- [79] A. Sasongko, I. N. Kumara, A. Wicaksana, F. Rousseau, O. Muller, Hardware context switch-based cryptographic accelerator for handling multiple streams, *ACM Transactions on Reconfigurable Technology and Systems (TRETS)* 14 (3) (2021) 1–25.
- [80] Nvidia jetson tx2 embedded systems, <https://www.xilinx.com/products/boards-and-kits/dk-v7-vc709-g.html>, accessed February 26, 2024.
- [81] A. A. Sützen, B. Duman, B. Şen, Benchmark analysis of jetson tx2, jetson nano and raspberry pi using deep-cnn, in: 2020 International Congress

on Human-Computer Interaction, Optimization and Robotic Applications (HORA), IEEE, 2020, pp. 1–5.

- [82] T. Amert, N. Otterness, M. Yang, J. H. Anderson, F. D. Smith, Gpu scheduling on the nvidia tx2: Hidden details revealed, in: 2017 IEEE Real-Time Systems Symposium (RTSS), IEEE, 2017, pp. 104–115.
- [83] Nvidia xavier series system-on-chip, <https://www.nvidia.com/en-us/autonomous-machines/embedded-systems/jetson-xavier-series>, accessed August 14, 2024.
- [84] Nvidia orin series system-on-chip, <https://www.nvidia.com/en-us/autonomous-machines/embedded-systems/jetson-orin>, accessed August 14, 2024.
- [85] Nvidia cuda c++ programming guide, [https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#compute-modes\(v12.6\)](https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#compute-modes(v12.6)), accessed September 16, 2024.