# High-Performance ARM-on-ARM Virtualization for Multicore SystemC-TLM-Based Virtual Platforms

Nils Bosbach* , Rebecca Pelke* , Niko Zurstraßen* , Jan Henrik Weinstock† ,
Lukas Jünger† , Rainer Leupers*
*RWTH Aachen University, Aachen, Germany
†MachineWare GmbH, Aachen, Germany
*{bosbach, pelke, zurstrassen, leupers}@ice.rwth-aachen.de     †{jan, lukas}@mwa.re

*Abstract*—The increasing complexity of hardware and software requires advanced development and test methodologies for modern systems on chips. This paper presents a novel approach to ARM-on-ARM virtualization within SystemC-based simulators using Linux's KVM to achieve high-performance simulation.

By running target software natively on ARM-based hosts with hardware-based virtualization extensions, our method eliminates the need for instruction-set simulators, which significantly improves performance. We present a multicore SystemC-TLM-based CPU model that can be used as a drop-in replacement for an instruction-set simulator. It places no special requirements on the host system, making it compatible with various environments.

Benchmark results show that our ARM-on-ARM-based virtual platform achieves up to 10 x speedup over traditional instruction-set-simulator-based models on compute-intensive workloads. Depending on the benchmark, speedups increase to more than 100 x.

*Index Terms*—**ARM-on-ARM, KVM, SystemC, TLM**

## I. INTRODUCTION

Driven by Moore's Law, the complexity of Hardware (HW) and Software (SW) has steadily increased over the past decades. This trend requires continuous evolution in the development and testing processes for modern Systems-on-Chips (SoCs). A technology that significantly supports companies in their SW development is virtual prototyping.

Virtual prototyping involves creating a simulator of the entire system, referred to as a Virtual Platform (VP) or Full-System Simulator (FSS) that behaves like the SoC under design. This VP can serve as a target platform for SW development while the HW is still being designed. One example is the pre-silicon development of device drivers using a virtual model of the device within a VP. VPs offer several other advantages including unlimited scalability, deep introspection, insightful tracing facilities, and seamless integration into automated Continuous Integration/Continuous Delivery (CI/CD) workflows.

One critical aspect of a VP's usability is its performance. High performance is essential for real-time debugging and automated testing. Typically, one of the most compute-intensive models within a VP is the CPU model executing the target SW. When the host[1]-machine architecture differs from that of the target[2] system, an Instruction-Set Simulator (ISS) within the CPU model translates instructions from the target Instruction-Set Architecture (ISA) to the host ISA.

---
[1]host: architecture that runs the simulation
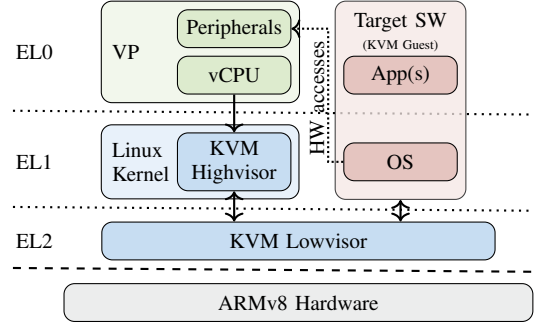[2]target: architecture to be simulated



Fig. 1: vCPU integration using KVM.

However, if both architectures match and Virtualization Host Extensions (VHE), such as ARMv8's VHE [1], are available on the host, the target SW can be executed natively without the need for an ISS. Figure 1 illustrates how Linux's hypervisor Kernel Virtual Machine (KVM) can be used by the virtual CPU (vCPU) to run the target SW natively on the host. Accesses to memory-mapped peripherals are trapped by KVM and forwarded to the CPU model running in the VP process on the host. This ensures that the target SW accesses the virtual HW instead of the physical host HW.

The industry-standard framework for building VPs is SystemC [2]. SystemC is a C++ library providing standardized interfaces that are crucial for compatibility across different simulations and tools. Transaction-Level Modeling (TLM) further extends the standard by providing abstract interfaces to model communication between different modules.

Building on previous work that integrated KVM into a SystemC-based CPU model [3], we add multicore support and eliminate limitations this approach has. With Apple's transition from Intel-based x86 to ARM-based Apple Silicon processors in 2020 [4], ARM-based HW that can be used as a platform for ARM-on-ARM (AoA)-based VPs has become widely available.

In this paper, we present:
- A multicore, SystemC-TLM-based CPU model that uses KVM and runs the simulated cores in parallel on the host. It serves as a drop-in replacement for an ISS without requiring any other adjustments to the VP.
- An implementation independent of performance counters.
- An alternate approach to custom kernel patches enhancing simulation performance.
- Benchmark results demonstrating our AoA-based VP's effectiveness compared to traditional ISSs.

arXiv:2505.12987v2 [cs.SE] 24 Jun 2025

## II. BACKGROUND & RELATED WORK

The field of virtual prototyping for embedded SW development has seen significant advancements in recent years, particularly in the context of ARM-based virtualization and SystemC-based simulations. This section provides an overview of existing literature and research in this domain, highlighting key contributions but also limitations of previous work.

### A. Virtual Prototyping for Software Development

Nowadays, the creation of VPs during the design process of an SoC has become essential. VPs model the behavior of the full system allowing unmodified target SW to be executed. This facilitates debugging and tracing [5] of the target SW.

SystemC has emerged as the industry-standard framework for building VPs due to its standardized interfaces and compatibility across different simulations and tools [2]. It provides basic module classes, a scheduler to simulate parallelism, and a concept of time. TLM extends SystemC by providing abstract interfaces to model communication between modules.

To further extend SystemC's features, modeling libraries are available that add frequently needed parts, components, and convenience functions. One example is the open-source Virtual Components Modeling Library (VCML) [6].

While SystemC-based simulations are widely used, alternative solutions, such as QEMU [7], exist. QEMU is a CPU-centric simulator supporting various host and target architectures. Compared to frameworks like SystemC, it lacks standardization which makes it challenging to integrate new models. Additionally, QEMU does not have a concept of simulation time which limits the detail level of the models. To overcome these limitations, solutions exist that wrap the processor model of QEMU in a SystemC module. An example of such a project is the open-source ARMv8 Virtual Platform (AVP64) [8], [9], which is used as a ISS-based reference system in this work.

### B. ARM-on-ARM Virtualization

Traditional CPU models use an ISS to translate instructions from the target ISA to the host ISA [7], [10]–[12]. Other approaches use recompilation of the target SW to the host ISA followed by native execution [13], [14]. While these approaches avoid the overhead of an ISS, they come with limitations such as required source-code access of the target SW.

A CPU model needs to execute the instructions of the target SW, simulate the CPU state including multiple Exception Levels (ELs), accept interrupts, and perform the address translations of the Memory Management Unit (MMU). Virtualization extensions allow to perform all these steps in HW on the host without requiring a complex SW solution. They add a layer of abstraction between the HW and the Operating System (OS). A hypervisor running in this layer can allow so-called guests to run side-by-side with the main OS (see Figure 1). Linux offers the hypervisor KVM [15] to use these virtualization features.

In 2020, the authors of [3] demonstrated how Linux's KVM can be used to build a SystemC-based CPU model that executes the target SW natively on an ARM host without needing an ISS. The work showed a basic proof of concept of a single-core VP. Their VP reached speedups of up to 2.57 x compared
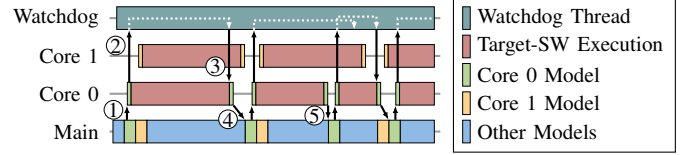


Fig. 2: Multithreaded approach using a SW-based watchdog.

to an ISS-based VP, although the ARM HW they used was less performant compared to the x86 machine running the ISS.

In summary, VPs offer numerous advantages SW development. High simulation performance is essential for efficient development and testing. Previous work showed that using Linux's KVM for SystemC-based CPU models results in huge speedups compared to traditional ISS-based models. However, the implementation relied on additional HW features and the possibility of applying custom kernel patches. In this work, we present how to remove these limitations and extend the CPU model to a multicore model.

## III. APPROACH

Our approach is sketched in Figure 2 for a dual-core setup. In a traditional SystemC-based simulation, only the *main thread* exists, where different models are simulated sequentially. For our parallelized approach, the SystemC kernel and all models except our CPU model run in the main thread. The CPU model has a placeholder in the main thread that is used for synchronization and the communication with other models. For the actual work, a separate thread is used. Details of this parallelization technique can be found in [16], [17].

For fast FSSs, the *loosely-timed* coding style is usually used to abstract timing for increased performance [2]. A technique used with this coding style is *temporal decoupling*. It permits SystemC processes to run ahead of the global simulation time before they need to synchronize again. A parameter called *quantum* defines how far a process is allowed to run ahead to control the trade-off between performance and accuracy. Temporal decoupling allows the CPU model to execute an entire quantum of instructions before synchronization is needed.

When the placeholder *SC_THREAD* of a CPU model is scheduled by SystemC, instructions of the target SW can be simulated. The worker running in the asynchronous thread, e.g., *Core 0* in Figure 2, is informed to execute a defined number of instructions ①. While other models continue their simulation in the main thread, a SW-based watchdog running in a different thread is programmed to stop the execution of target SW after the specified wall-clock-time interval ②. Then, the target SW is executed natively on the host using the virtualization features. When the watchdog expires ③, the execution of the target SW is suspended and the placeholder *SC_THREAD* running in the main thread is notified ④.

Depending on the instructions of the target SW, the execution might be suspended before watchdog expiration. This is, e.g., the case when the target SW needs to access a memory-mapped region of a peripheral. Then, the access to the peripheral is performed in the main thread ⑤. Afterwards, the execution of the target SW can be continued. In case of an early exit, the programmed watchdog is not needed (details in Section IV-B).

In this work, we present a technique to force early accesses when the Linux of the target SW executes an idle loop. The Linux idle loop uses the Wait For Interrupt (WFI) instruction to hint a possible suspension until an interrupt is signaled. We can use this forced early access to skip the execution of idle loops and thereby speed up the simulation. Details can be found in Section IV-C.

## IV. Implementation

We build our implementation of the KVM-based processor model on top of the SystemC-TLM-based, open-source VCML [6] project. VCML offers the `processor` class as a starting point to integrate an instruction interpreter, simulator, or executor into a loosely-timed SystemC module. It implements a basic simulation loop that calls the `virtual void simulate(size_t cycles)` function to simulate the specified number of instructions. This function is used to call either the ISS or, in our case, KVM. Additionally, parallel execution of the simulate function can be activated to offload the execution to a separate thread on the host machine [16], [17].

Communication with KVM is mainly done through system calls. To efficiently handle memory accesses of the target SW, a memory region from the VP process, can be mapped to the KVM guest. We use this feature to map the virtual memory model of the VP to the guest environment. The region is queried using TLM-Direct Memory Interface (DMI), similar to how an ISS would operate. This allows the native execution of load and store instructions to memory.

### A. Working Principle

A basic overview of the CPU model's simulate function is given in Figure 3. When parallel execution is activated, the call is executed in a separate thread for each simulated CPU core to allow other models of the simulation to run in parallel. Otherwise, the *CPU-Core Thread* does not exist, and the functions are directly executed in the SystemC thread.

First, the allowed runtime of KVM is calculated from the passed number of cycles. As usual for instruction-accurate simulators, a constant average execution time per instruction is assumed. The clock frequency is obtained from the virtual clock that is connected to the SystemC processor module. The watchdog is programmed to expire after the calculated amount of time (more details on the watchdog follow in Section IV-B).

Pending interrupts are injected. The (wall-clock) timestamp is stored in a local variable and the execution of KVM is started using the `KVM_RUN` call. KVM then executes the guest code. Certain events are trapped by KVM that lead to a return to user space and therefore an exit of the `KVM_RUN` call. The relevant events for this work are Memory-Mapped Input/Output (MMIO) accesses, breakpoint hits, and received Linux signals.

After an exit of `KVM_RUN`, the internal watchdog ID is incremented (more on this in Section IV-B), and the duration of the `KVM_RUN` call is determined by subtracting the current (wall-clock) timestamp from the previously stored one. This runtime is used as an approximation for the number of executed instructions and thereby cycles.
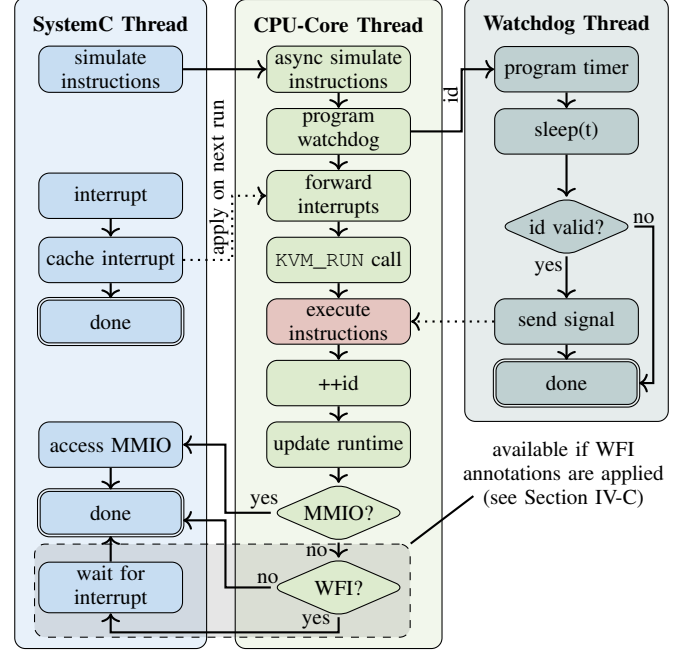


Fig. 3: vCPU model execution loop.

Then, the cause of the `KVM_RUN` exit is determined. In case of an MMIO access, a TLM transaction is sent to the corresponding peripheral. Since, according to the SystemC standard, interactions with other models need to be performed in the main thread, this access has to be shifted back in case of parallel execution [2], [16]. If the exit was caused by an idling core that executed a WFI instruction, performance optimization can be applied (more details in Section IV-C).

Once the MMIO handling is done, the WFI hint has been processed, or if the exit has been caused by the watchdog, the `simulate` function returns. The loop of the VCML processor class will continue and call the `simulate` function again.

### B. Software-Based Watchdog Timer

To stop the execution after the end of a quantum, a SW-based watchdog timer is used (see Section III). The vCPU clock is used to convert the quantum into the number of instructions that can be executed before synchronizing. Once the number of instructions has been executed, a Linux signal is sent to the thread that executes the `KVM_RUN` call to stop the execution.

Recent approaches have used Linux's CPU-performance-counter-API *perf* to count the executed instructions in guest mode and send a signal once a threshold is exceeded [3]. While this method provides high accuracy, it depends on specific HW features that may not be universally available. For instance, Apple's CPUs contain custom Performance Monitoring Units (PMUs) instead of the standard ARM one. In Asahi Linux [18], perf cannot be configured to send a signal once a specified amount of guest instructions has been executed. Therefore, we use an approach that does not rely on HW-specific features.

Before invoking `KVM_RUN`, we set up the SW-based watchdog timer. The watchdog timer is shared between all cores and runs in a separate thread (see Figures 2 and 3). It calls the `kick` function after a predefined timeout that is calculated by dividing

the maximum number of instructions to be executed by the vCPU clock frequency. The kick function and the scheduling of the timer are shown in Listing 1.

```
1  // kick KVM
2  void cpu::kick(unsigned int id) {
3      if (id == m_kickid)
4          pthread_kill(m_self, SIGUSR1);
5  }
6  // schedule watchdog
7  watchdog(timeout, [&, id = m_kickid]() -> void {
8              kick(id); });
```

Listing 1: Kick function that is called by the watchdog.

Our CPU class maintains an internal ID counter (`m_kickid`) that is incremented after each `KVM_RUN` to identify a run. When scheduling the watchdog timer, we pass this current ID value (Lines 7 and 8). Upon expiration, the ID is used to only send the `SIGUSR1` signal to exit KVM if the corresponding `KVM_RUN` is still active (Lines 3 and 4). This approach effectively limits the maximum KVM run time.

By employing this method, we achieve robust synchronization without relying on hardware-specific performance counters. Consequently, our solution is versatile and compatible with various environments, such as Asahi Linux on Apple Silicon.

### C. WFI Annotations

The ARM architecture includes a WFI instruction, that is commonly used in idle loops within OSs to suspend CPU activity until an interrupt occurs [1]. In a SystemC-based, event-driven simulation, pausing the execution of the CPU model during idle periods and resuming it upon receiving an interrupt can significantly enhance simulation performance by skipping idle time instead of simulating it.

However, KVM does not inherently notify the process that started a KVM guest about executed WFI instructions. The WFI instructions are trapped by KVM and Linux's scheduler is called to schedule another user thread. Therefore, modifications in KVM or workarounds are necessary to get notified of WFI instruction execution to then use SystemC's feature to pause the model until the next interrupt occurs.

Previous approaches used custom patches to send notifications to user space [3]. While this is an effective approach for research purposes, this method comes with several drawbacks:

- *Security concerns*: Stringent security policies may prevent companies from applying custom patches.
- *Maintenance overhead*: Custom patches need reapplication with every kernel update, complicating maintenance.
- *Mainline-integration challenges*: Efforts to integrate such patches into the mainline Linux kernel failed.

To address these limitations, we propose an alternative way to trap WFI instructions without requiring kernel patches. This approach works when a Linux OS is executed on the target system since Linux uses the WFI instruction only in its idle loop. Our solution involves searching for specific symbols within the target SW's Executable and Linkable Format (ELF) file during VP startup and setting breakpoints accordingly. The steps are as follows:
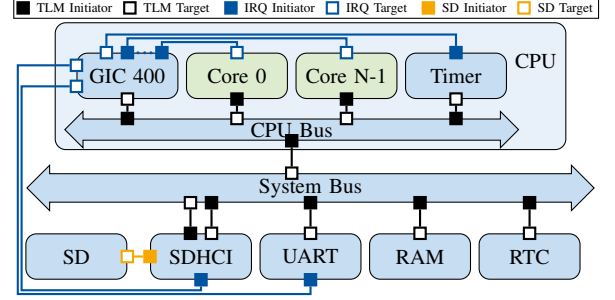


Fig. 4: AoA-based VP.

1) *Symbol search*: Identify the `cpu_do_idle` symbol within the target-SW ELF file.
2) *Breakpoint setting*: Locate the WFI instruction inside this function and set a breakpoint on it.
3) *Instruction check*: When a vCPU reaches this breakpoint, KVM exits back to user space.
4) *Program counter verification*: Verify if the program counter matches the address of the WFI instruction to distinguish it from other breakpoints set by users.

Upon confirming that a CPU intends to execute a WFI instruction, we use SystemC's features to suspend model execution until an interrupt is signaled. This technique avoids the simulation of idle loops, thereby improving the overall simulation efficiency and performance.

We refer to this method as *WFI annotation*. It provides a robust way to handle idle states in multicore simulations without compromising security or maintainability.

## V. RESULTS

To evaluate the performance results of our approach, we integrated the KVM-based CPU model into a VP. Figure 4 shows an overview of the VP architecture. Each green-colored CPU core launches a KVM-based guest to execute target instructions. The vCPU consists of 1 to 8 KVM-based cores, a Generic Interrupt Controller (GIC) 400, and a memory-mapped timer. The simulated peripherals include a SD Host Controller Interface (SDHCI) device with a virtual Secure Digital (SD) card, an Universal Asynchronous Receiver/Transmitter (UART) interface for user interaction, Random-Access Memory (RAM), and a Real-Time Clock (RTC). All peripherals are taken from the open-source VCML library [6]. Communication between the models is realized via TLM sockets and protocols for memory accesses, interrupts, and SD-card operations.

We conducted our experiments on a 10-core *Apple Mac mini* equipped with an *M2 Pro* processor, 16 GB RAM, and 512 GB Solid State Drive (SSD) memory. The M2 Pro processor has 6 high-performance *Avalanche* (3.7 GHz) cores and 4 efficiency *Blizzard* (3.4 GHz) cores. For comparison against traditional ISS-based CPU models, we used the open-source AVP64 [8], [9] running on an *AMD Ryzen 9 3900X* (3.8 GHz, 4.6 GHz boost) with 12 cores and 24 HW threads. Both VPs support a parallel execution mode that runs the CPU cores in separate threads [16], [17]. Different quantum values are used to steer the temporal decoupling of the simulated models.
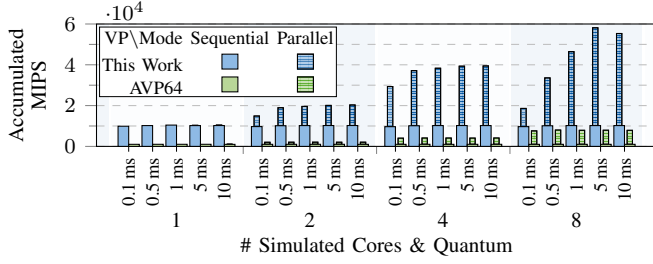
Fig. 5: Bare-metal Dhrystone [19] benchmark.

## A. Bare-Metal Dhrystone Benchmark

The first benchmark executed is a bare-metal Dhrystone [19], which is single threaded and integer based. For multicore systems, each core executes its own instance of Dhrystone. Thus, a parallel Dhrystone is an optimally parallelizable, compute-intensive workload that does not involve any communication.

Figure 5 shows the measured accumulated Million Instructions Per Second (MIPS) values for different core counts, quantum values, and parallelization settings. For a single-core VP, parallelization does not yield performance benefits as the CPU is the only compute-intensive model that runs in the main thread. Our AoA VP achieves nearly 10,000 MIPS, which is about 10 x the performance of AVP64.

For dual-core setups with parallel execution enabled, the performance effectively doubles due to simultaneous simulation of both cores. However, smaller quantum values lead to decreased AoA performance due to increased synchronization overheads involving EL-switching for KVM entries and exits [20]. This effect is even more visible in quad- and octa-core setups.

In octa-core configurations, limited host-machine performance cores (6 in total) reduce the achievable speedups since some simulated cores have to run on efficiency cores.
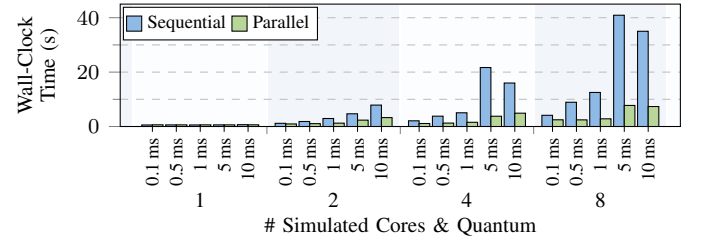
In summary, the parallel Dhrystone benchmark shows that for compute-intensive workloads, the AoA VP achieves up to 10 x speedup over the ISS-based AVP64. Parallelized execution of the simulated cores significantly enhances the performance. The optimal speedup can be reached for dual and quad-core setups. For octa-core setups, the limited number of performance cores of the host machine reduces the achievable speedup.
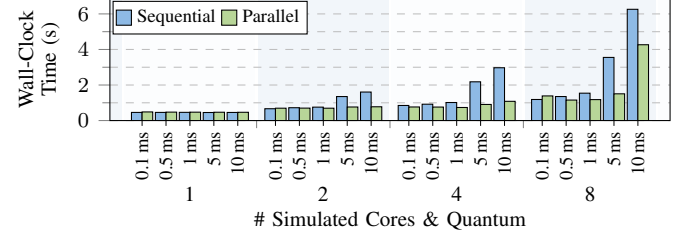
## B. Linux's Boot Process

The next benchmark we look at is the boot process of a Buildroot-based Linux [21]. In contrast to the parallel Dhrystone, this benchmark represents a mostly sequential workload where one core performs the boot tasks while others idle.

Figure 6a illustrates the wall-clock time required for different core counts without WFI annotations (see Section IV-C, WFI is handled by KVM). While the single-core VP boots Linux in approximately 0.6 s, multicore setups take considerably longer. A non-parallelized boot of an octa-core VP can take up to 40 s. The reason for this is that in addition to the core that performs the boot, the idle loops of the other cores need to be simulated. For larger quantum values, synchronization between the cores is more complicated which leads to increased runtime. This has also been observed by previous works [22]–[24].

When parallelization is activated, the idling cores can be simulated in parallel which reduces the needed amount of wall-



(a) Without WFI annotation.



(b) With WFI annotations.

Fig. 6: Buildroot Linux boot durations for AoA.

clock time. However, booting a multicore system is still slower than booting a single-core one due to the additional overhead.

Apart from applying parallelization, another way to optimize the idle-loop behavior is to annotate the WFI instruction (see Section IV-C). Figure 6b shows how WFI annotations effectively increase the performance of the Linux boot. Instead of simulating the idle loop, the simulation of the core is suspended until an interrupt is signaled. This reduces the time needed for a Linux boot to under 1 s for dual and quad-core setups. For the octa-core VP, speedups achieved by WFI annotation range from 1.78 x for the 100 μs parallel version up to 11.5 x for the 5 ms sequential version.

The Linux-boot benchmark shows that the simulation of idle loops drastically reduces the performance of the sequential simulation. Parallelizing the simulation helps to counteract this effect. However, the best results are achieved when idle loops are not simulated but annotated, so the simulation of the core can be suspended. When both techniques are combined, the needed time for a multicore Linux boot can be kept small.

## C. User-Space Benchmarks

After analyzing fully parallelizable and a predominantly sequential workloads, we now look at common user-space benchmarks executed within a Linux environment. The results are depicted in Figure 7 for a 1 ms quantum with parallel execution enabled, and various core configurations. All listed benchmarks are executed on AVP64 and our AoA VP. The wall-clock time it takes to execute the benchmarks is measured. On the y-axis of Figure 7, one can see the speedup of AoA compared to the execution on a similarly configured AVP64 (quantum, number of cores, parallelized execution).

In addition to the user space benchmarks, we also include results for the bare-metal Dhrystone and Linux-boot process. As observed in Section V-A, the Dhrystone benchmark shows a dip in speedup for eight simulated cores due to the number of available performance cores. For the Linux boot process, increased core counts reduce the speedup because trapping WFI
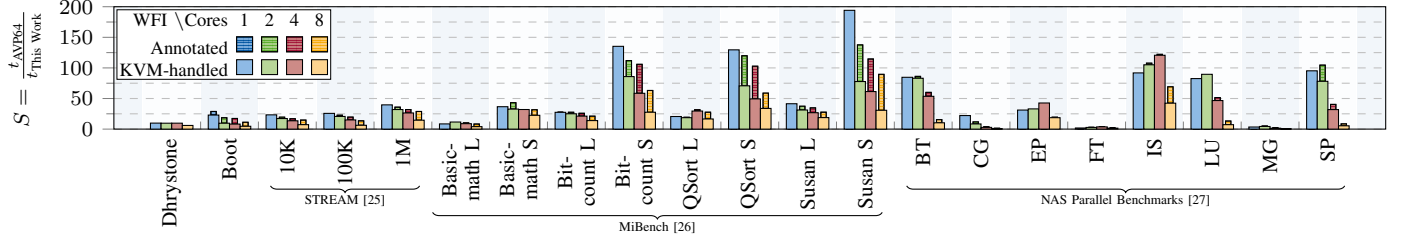
Fig. 7: Benchmark speedup $S$ of AoA compared to AVP64 for 1 ms quantum and activated parallel execution.

instructions is less expensive for ISSs than for AoA due to the needed context switching for exiting KVM [20].

*1) STREAM:* The first user space benchmark tested is *STREAM* [25], which measures memory bandwidth by executing numerous load and store instructions [28]. We run the benchmark with different array sizes ranging from 10,000 (*10K*) to 1,000,000 (*1M*) elements. For ISS-based simulators like AVP64, these instructions require SW-based MMU translations from virtual to physical addresses. This task incurs significant overhead. In contrast, AoA models leverage the HW-provided two-stage address translation process of the host MMU that handles these translations natively [29], [30].

*2) MiBench Suite:* Next, we evaluate benchmarks from the *MiBench* [26] suite, representing automotive and industrial control workloads. Theses single-threaded benchmarks do not benefit from a multicore setup. Therefore, annotating WFI instructions in the idle loop is essential to limit the performance loss on multicore setups. Speedups compared to the ISS-based AVP64 range significantly from approximately 8 x for *Basicmath L* up to 165x for *Susan S* on single-core VPs.

MiBench provides both large (*L*) and small (*S*) versions of each benchmark type. The variants only differ in size of the input but not in the task itself. It can be observed that the smaller variants achieve higher speedups than the larger ones. This indicates that the huge speedup obtained by the small variants is not mainly caused by the computation itself. For the Dynamic Binary Translation (DBT)-based ISS used by AVP64, basic blocks of the target SW are translated to the host ISA. Once they have been translated, they are cached, so the translation is only needed once. When the size of the input for a workload is increased, usually the loops of the algorithm are just executed more often. That means the translation overhead of the ISS has a larger impact when the translated blocks are executed less often. The ISS therefore gets more efficient for increased workload-input sizes. For AoA, those effects do not exist whereby the speedup of AoA compared to AVP64 is reduced when the ISS gets more efficient.

*3) NAS Parallel Benchmarks:* Lastly, we examine multicore benchmarks from the NAS Parallel Benchmarks (NPB) [27] suite. These benchmarks distribute their workloads across multiple threads using Open Multi-Processing (OpenMP) [31]. Since the benchmarks use all cores of the simulated system, the system does not spend much time in idle mode. Therefore, WFI annotation does not significantly improve the performance.

Workloads that involve a lot of synchronization and therefore communication between the cores, like the *CG*, *FT*, and *MG* benchmarks, cause more overhead than the other workloads.

However, with a minimum speedup of 1.8 x for the FT benchmark, AoA is still faster than AVP64.

In summary, the benchmark results show that AoA is a promising approach that results in huge speedups compared to traditional ISS-based VPs. The achieved speedup of 2.57 x of [3] for a bare-metal Coremark compared to an ISS-based VP can be increased to up to 10 x for Dhrystone on modern ARM HW with parallel execution. Apart from the observed speedups, the results show that idle-loop annotations or WFI trapping are essential to limit the performance loss of single-threaded workloads on a multicore VP.

## VI. CONCLUSION & FUTURE WORK

In this paper, we presented a novel approach to AoA virtualization within SystemC-based VPs, leveraging Linux's KVM to enhance performance. By running the target SW natively on ARM-based hosts with VHE, we eliminate the need for an ISSs and thereby significantly improve the simulation performance. Our SystemC-TLM-based multicore CPU model operates independently of host performance counters or custom kernel patches, which makes it more independent of the environment.

Benchmark results demonstrated that our AoA-based VP approach achieves substantial speedups compared to traditional ISS-based models. For compute-intensive workloads, such as parallel Dhrystone, our solution reaches up to 10 x the performance over an ISS-based VP. Additionally, by annotating WFI instructions, we were able to further optimize the simulation of idle loops during processes like Linux booting, achieving significant reductions in required wall-clock time.

Overall, our work shows that AoA virtualization is a promising technique for SW development and testing. We presented that the annotation of idle loops improves the simulation performance of certain workloads. While our annotation-based approach works well for Linux-based target SW, it needs to be adapted for other workloads. For future work, we plan to further improve or automate this way of trapping WFI instructions to be able to handle WFI instructions without manual annotations.

Another limitation of our approach is that KVM currently only supports EL0 and EL1 for the guest. This is sufficient for most workloads including OSs and user-space SW. However, hypervisors running in EL2 and ARM's trusted firmware running in EL3 cannot be executed with our approach.

In addition, instruction emulation can be added to support new instructions that may emerge in the future that are not supported by the host. Furthermore, the approach can be extended to other architectures that have a virtualization extension, such as RISC-V-on-RISC-V simulation.

REFERENCES

[1] ARM Ltd, "Arm Architecture Reference Manual DDI 0487K.a," Mar. 2024.

[2] IEEE Standards Association and others, "IEEE Standard for Standard SystemC® Language Reference Manual," *IEEE Std 1666-2023 (Revision of IEEE Std 1666-2011)*, pp. 1–618, Sep. 2023, conference Name: IEEE Std 1666-2023 (Revision of IEEE Std 1666-2011). [Online]. Available: https://doi.org/10.1109/IEEESTD.2023.10246125

[3] L. Junger, J. L. Malte Bolke, S. Tobies, R. Leupers, and A. Hoffmann, "ARM-on-ARM: Leveraging Virtualization Extensions for Fast Virtual Platforms," in *2020 Design, Automation & Test in Europe Conference & Exhibition (DATE)*. Grenoble, France: IEEE, Mar. 2020, pp. 1508–1513. [Online]. Available: https://doi.org/10.23919/DATE48585.2020.9116573

[4] R. Iyengar, "Apple details new MacBook Air, MacBook Pro and Mac Mini – all powered by in-house silicon chips | CNN Business," Nov. 2020. [Online]. Available: https://www.cnn.com/2020/11/10/tech/apple-silicon-chips-mac/index.html

[5] N. Bosbach, L. Jünger, J. M. Joseph, and R. Leupers, "NISTT: A Non-Intrusive SystemC-TLM 2.0 Tracing Tool," in *2022 IFIP/IEEE 30th International Conference on Very Large Scale Integration (VLSI-SoC)*, Oct. 2022, pp. 1–6, iSSN: 2324-8440. [Online]. Available: https://doi.org/10.1109/VLSI-SoC54400.2022.9939578

[6] MachineWare, "machineware-gmbh/vcml," Mar. 2024, original-date: 2018-01-22T10:24:21Z. [Online]. Available: https://github.com/machineware-gmbh/vcml

[7] F. Bellard, "QEMU, a fast and portable dynamic translator." in *USENIX annual technical conference, FREENIX Track*, vol. 41. California, USA, 2005, pp. 10–5555, issue: 46.

[8] L. Jünger, J. H. Weinstock, R. Leupers, and G. Ascheid, "Fast SystemC Processor Models with Unicorn," in *Proceedings of the Rapid Simulation and Performance Evaluation: Methods and Tools*, ser. RAPIDO '19. New York, NY, USA: Association for Computing Machinery, Jan. 2019, pp. 1–6. [Online]. Available: https://doi.org/10.1145/3300189.3300191

[9] L. Jünger, "An ARMv8 Virtual Platform (AVP64)," May 2023, original-date: 2020-04-09T15:34:12Z. [Online]. Available: https://github.com/aut0/avp64

[10] L. Jünger, J. H. Weinstock, and R. Leupers, "SIM-V: Fast, Parallel RISC-V Simulation for Rapid Software Verification," in *Proceedings of DVCon Europe 2022*, Munich, 2022. [Online]. Available: https://dvcon-proceedings.org/document/sim-v-fast-parallel-risc-v-simulation-for-rapid-software-verification/

[11] P. Magnusson *et al.*, "Simics: A full system simulation platform," *Computer*, vol. 35, no. 2, pp. 50–58, Feb. 2002, conference Name: Computer. [Online]. Available: https://doi.org/10.1109/2.982916

[12] N. Binkert *et al.*, "The gem5 simulator," *SIGARCH Comput. Archit. News*, vol. 39, no. 2, pp. 1–7, Aug. 2011. [Online]. Available: https://doi.org/10.1145/2024716.2024718

[13] A. Gerstlauer, "Host-compiled simulation of multi-core platforms," in *Proceedings of 2010 21st IEEE International Symposium on Rapid System Protyping*, Jun. 2010, pp. 1–6, iSSN: 2150-5519. [Online]. Available: https://doi.org/10.1109/RSP.2010.5656355

[14] H. Shen, M.-M. Hamayun, and F. Petrot, "Native Simulation of MPSoC Using Hardware-Assisted Virtualization," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 31, no. 7, pp. 1074–1087, Jul. 2012, conference Name: IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems. [Online]. Available: https://doi.org/10.1109/TCAD.2012.2187526

[15] C. Dall and J. Nieh, "KVM/ARM: the design and implementation of the linux ARM hypervisor," *SIGPLAN Not.*, vol. 49, no. 4, pp. 333–348, Feb. 2014. [Online]. Available: https://doi.org/10.1145/2644865.2541946

[16] N. Bosbach, N. Zurstraßen, R. Pelke, L. Jünger, J. H. Weinstock, and R. Leupers, "Towards High-Performance Virtual Platforms: A Parallelization Strategy for SystemC TLM-2.0 CPU Models," in *Design Automation Conference*, Jun. 2024. [Online]. Available: https://doi.org/10.1145/3649329.3658257

[17] N. Bosbach, R. Pelke, N. Zurstraßen, L. Junger, J. H. Weinstock, and R. Leupers, "Work-in-Progress: A Generic Non-Intrusive Parallelization Approach for SystemC TLM-2.0-based Virtual Platforms," in *Proceedings of the 2023 International Conference on Hardware/Software Codesign and System Synthesis*. Hamburg Germany: ACM, Sep. 2023, pp. 42–43. [Online]. Available: https://doi.org/10.1145/3607888.3608596

[18] Asahi Linux, "Asahi Linux." [Online]. Available: https://asahilinux.org/

[19] R. P. Weicker, "Dhrystone benchmark: rationale for version 2 and measurement rules," *AcM SIGPLAn notices*, vol. 23, no. 8, pp. 49–62, 1988, publisher: ACM New York, NY, USA.

[20] J. Engblom, "How the Intel® Simics® Simulator Executes Instructions," Nov. 2023, section: Software. [Online]. Available: https://community.intel.com/t5/Blogs/Products-and-Solutions/Software/How-the-Intel-Simics-Simulator-Executes-Instructions/post/1543049

[21] Buildroot, "Making Embedded Linux Easy," 2024. [Online]. Available: https://buildroot.org/

[22] N. Zurstraßen, R. Brandhofer, J. Cubero-Cascante, N. Bosbach, L. Jünger, and R. Leupers, "The Optimal Quantum of Temporal Decoupling," in *2024 29th Asia and South Pacific Design Automation Conference (ASP-DAC)*, 2024, pp. 686–691. [Online]. Available: https://doi.org/10.1109/ASP-DAC58780.2024.10473967

[23] N. Zurstraßen, R. Brandhofer, J. Cubero-Cascante, N. Bosbach, L. Jünger, and R. Leupers, "The art of temporal decoupling," *Integration*, vol. 101, p. 102314, Mar. 2025. [Online]. Available: https://doi.org/10.1016/j.vlsi.2024.102314

[24] J. Engblom, "Some Notes on Temporal Decoupling (Reposted) – Observations from Uppsala," Mar. 2022. [Online]. Available: https://jakob.engbloms.se/archives/3467

[25] J. D. McCalpin, "Memory Bandwidth and Machine Balance in Current High Performance Computers," *IEEE Computer Society Technical Committee on Computer Architecture (TCCA) Newsletter*, pp. 19–25, Dec. 1995.

[26] M. Guthaus, J. Ringenberg, D. Ernst, T. Austin, T. Mudge, and R. Brown, "MiBench: A free, commercially representative embedded benchmark suite," in *Proceedings of the Fourth Annual IEEE International Workshop on Workload Characterization. WWC-4 (Cat. No.01EX538)*. Austin, TX, USA: IEEE, 2001, pp. 3–14. [Online]. Available: https://doi.org/10.1109/WWC.2001.990739

[27] NASA Jet Propulsion Laboratory, "nasa-jpl/embedded-gcov," Oct. 2023, original-date: 2022-02-02T19:25:26Z. [Online]. Available: https://github.com/nasa-jpl/embedded-gcov

[28] N. Bosbach, L. Jünger, R. Pelke, N. Zurstraßen, and R. Leupers, "Entropy-Based Analysis of Benchmarks for Instruction Set Simulators," in *Proceedings of the DroneSE and RAPIDO: System Engineering for constrained embedded systems*, ser. RAPIDO '23. New York, NY, USA: Association for Computing Machinery, Apr. 2023, pp. 54–59. [Online]. Available: https://doi.org/10.1145/3579170.3579267

[29] ARM, "Arm System Memory Management Unit Architecture Specification," no. 3, Feb. 2024.

[30] ARM Ltd, "Armv8-A virtualization 102142," 2019.

[31] R. Chandra, L. Dagum, D. Kohr, R. Menon, D. Maydan, and J. McDonald, *Parallel programming in OpenMP*. Morgan kaufmann, 2001.