

Pac-Sim: Simulation of Multi-threaded Workloads using Intelligent, Live Sampling

Changxi Liu[†]
National University of Singapore

Alen Sabu[†]
National University of Singapore

Akanksha Chaudhari
University of Wisconsin-Madison

Qingxuan Kang
National University of Singapore

Trevor E. Carlson
National University of Singapore

Abstract—High-performance, multi-core processors are the key to accelerating workloads in several application domains. To continue to scale performance at the limit of Moore’s Law and Dennard scaling, software and hardware designers have turned to dynamic solutions that adapt to the needs of applications in a transparent, automatic way. For example, modern hardware improves its performance and power efficiency by changing the hardware configuration, like the frequency and voltage of cores, according to a number of parameters such as the technology used, the workload running, etc. With this level of dynamism, it is essential to simulate next-generation multi-core processors in a way that can both respond to system changes and accurately determine system performance metrics. Currently, no sampled simulation platform can achieve these goals of dynamic, fast, and accurate simulation of multi-threaded workloads.

In this work, we propose a solution that allows for fast, accurate simulation in the presence of both hardware and software dynamism. To accomplish this goal, we present *Pac-Sim*, a novel sampled simulation methodology for fast, accurate sampled simulation that requires no upfront analysis of the workload. With our proposed methodology, it is now possible to simulate long-running dynamically scheduled multi-threaded programs with significant simulation speedups even in the presence of dynamic hardware events. We evaluate *Pac-Sim* using the multi-threaded SPEC CPU2017, NPB, and PARSEC benchmarks with both static and dynamic thread scheduling. The experimental results show that *Pac-Sim* achieves a very low sampling error of 1.63% and 3.81% on average for statically and dynamically scheduled benchmarks, respectively. *Pac-Sim* also demonstrates significant simulation speedups as high as $523.5\times$ ($210.3\times$ on average) for the train input set of SPEC CPU2017.

I. INTRODUCTION

Computer architecture research heavily relies on simulations for design space exploration. However, microarchitectural simulation can become extremely time-consuming, particularly as the complexity of modern architectures has increased over time. This is especially true in the post-Dennard era, where architectures are rapidly evolving to incorporate complex dynamic optimization techniques at both the hardware and software levels to improve system performance gains at runtime. Hardware-based dynamic techniques such as dynamic cache reconfiguration [4], [51], [52], DVFS [33], [42], [45], TurboBoost [19] and power management [11] techniques trigger optimizations based on dynamically identified hardware

states (such as core frequency, cache reuse distance, etc.) to improve both energy-efficiency and overall performance of the system. Similarly, runtime information at the software level can be used to dynamically optimize code execution, to further enhance the system performance. Some of the recent efforts on software-based optimization focus on dynamically scheduling tasks among threads [27], [29] to ensure efficient resource utilization and employing just-in-time (JIT) compilation techniques [22], [46], [67], [70] that generate high-performance instructions to optimize program execution online. However, since these techniques utilize dynamic system state information in order to deploy optimizations at runtime, the execution behavior of an application (and, therefore, its performance) may vary vastly across multiple executions. Such variability can make it extremely difficult to determine the performance of a given workload using existing simulation methodologies.

Conventionally, sampled simulation has served as a reliable and efficient technique to accelerate the performance estimation of multi-threaded workloads. In order to achieve these results, most prior works relied on either (i) profile-driven sampling [17], [60], [64] or (ii) statistical sampling [43], [69]. Profile-driven sampled simulation methodologies such as SimPoint [64], BarrierPoint [17], and LoopPoint [60] split the execution of an application into a series of repeatable regions and cluster them based on their execution features. A representative element from each cluster is then analyzed or simulated in order to extrapolate the performance of the entire application. However, these methodologies incur a significant cost in terms of the preprocessing effort that is needed to identify representative regions. These costs include the time required to profile and cluster the execution features of all application regions, along with the storage required. While it has been previously argued that these costs are a one-time investment and will be amortized over multiple runs, this argument does not necessarily hold for systems that optimize code execution dynamically. In such cases, the program execution paths followed by an application may vary considerably due to changes of hardware and software parameters that are being optimized. Therefore, the profiling information collected for one specific run would not necessarily extend to the program execution paths followed in the subsequent runs.

On the other hand, sampled simulation methodologies such

[†]Changxi Liu and Alen Sabu contributed equally to this work.

as SMARTS [69] and PCantorSim [43] rely on statistical sampling techniques to speed up simulation-based performance measurements while meeting a given error bound. Unlike profile-driven sampling, these methodologies require minimal preprocessing and do not rely on the reproducibility of program execution paths. They are thus applicable to dynamically optimized systems. However, the simulation speedups achieved using these techniques are considerably lower than the profile-driven counterparts, and adjusting settings to achieve higher performance could lead to high errors.

For the above-mentioned reasons, it becomes challenging to sample and simulate generic multithreaded applications for dynamic hardware and software using existing methodologies. Architects need a simulation methodology that can dynamically adapt to changes in the system at runtime while accurately estimating the application’s performance without relying on the reproducibility of its execution. To this end, we propose Pac-Sim, a novel sampled simulation methodology that can, at runtime, efficiently analyze and sample the application to select the representative regions to be simulated in detail. The result is a methodology that enables both fast and accurate performance evaluation without the need for up-front analysis. We accomplish this by making use of an intelligent online¹ predictor and classifier that quickly and accurately decides whether the upcoming region needs to be simulated in detail.

In short, we make the following contributions:

- i. We propose Pac-Sim, a methodology that goes beyond prior sampled simulation techniques to be the first to allow for dynamic hardware and software support. The methodology requires no upfront analysis and relies on an online predictor for sampling decisions enabling the fast analysis of co-designed workloads. We will open-source the simulation framework for Pac-Sim upon acceptance.
- ii. We experimentally demonstrate that Pac-Sim consistently improves performance in terms of speedup and accuracy over prior works that use offline profiling. We also quantify the performance benefits obtained by Pac-Sim, showing that they outweigh its online analysis overheads.
- iii. We provide an extensive evaluation of Pac-Sim using standard benchmarks to compare against prior works and demonstrate best-in-class accuracy (average error of 1.63%). For the SPEC CPU2017 benchmarks (train inputs) running eight threads, we show a maximum serial speedup of $123.32\times$ ($26.09\times$ on average) and a maximum parallel speedup of $523.5\times$ ($210.3\times$ on average).
- iv. Finally, we showcase several case studies demonstrating that Pac-Sim is applicable to a number of research scenarios, including (but not limited to) the investigation of optimization techniques such as dynamically scheduled software, and improving research into dynamic hardware and hardware-software co-design.

The rest of the paper is organized as follows. In Section II, we discuss the relevant background and the challenges in-

involved in the simulation of dynamic applications on modern architectures. Section III presents the Pac-Sim methodology in detail. We then describe the experimental infrastructure in Section IV, followed by an extensive evaluation of Pac-Sim in Section V along with case studies to demonstrate the applicability of the proposed methodology. Finally, we conclude the paper in Section VII.

II. SIMULATING MODERN ARCHITECTURES

In this section, we provide the necessary background of sampled simulation. We also discuss the challenges in simulating modern workloads and how the existing sampling methodologies are insufficient to address them.

Sampling Single-threaded Workloads. Sampling and workload reduction techniques are extensively utilized in computer architecture research for the purpose of program characterization and to reduce simulation time. Sampling methodologies allow for the evaluation of a subset of the workload (a representative sample) in detail that can be used to reconstruct the performance of the whole workload accurately. These methodologies split the workload into different regions (or slices) based on predetermined conditions in order to identify a representative sample. Prior works that explored CPU workload sampling, like SimPoint [64] and SMARTS [69], tend to utilize fixed instruction counts to determine regions. However, instruction count-based techniques could lead to inconsistent and, therefore, invalid regions [2], [3], [16]. Some previous works [48], [71] proposed software phase markers that identify procedure and loop boundaries that correlate with phase changes to mark region boundaries instead of using fixed-sized regions. On the other hand, statistical sampling methods, including SimFlex [68], were proposed for multi-processor throughput workloads and use instruction count for statistical sampling, but these works do not appear to be generally extensible to synchronizing multi-threaded workloads [16]. In the presence of synchronizing threads, the application performance tends to vary more frequently, and the statistical confidences assuming Gaussian performance distribution, as shown in SMARTS or SimFlex, may not be applicable there [21].

Sampling Multi-threaded Workloads. Time-based sampling methodologies [5], [16] are the first to address the problem of sampling synchronizing multi-threaded applications. These methodologies, however, are slow, and as a result, they are not practical for handling realistic workloads. On the other hand, methodologies like BarrierPoint [17], TaskPoint [38], and LoopPoint [60] select specific program constructs, such as barrier synchronization primitives, task instances, and loops, respectively, to identify periodic behavior. This enables the utilization of representative-sized regions for simulation, regardless of the program’s length.

Feature Vectors. Profiling captures feature vectors up-front to characterize the execution behavior of an application across regions. Previous works have introduced several microarchitecture-independent feature vectors, of which basic block vectors (BBVs) [59], [64] are the most widely used for

¹We use the terms “online” and “offline” to distinguish between events that occur during and prior to the simulation of an application, respectively.

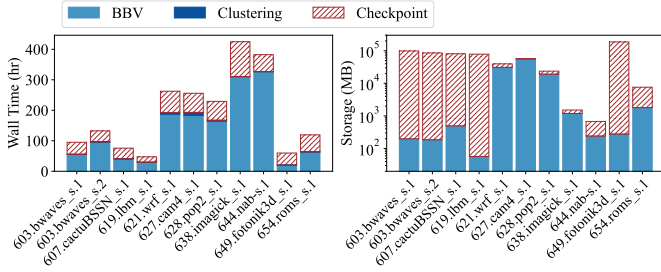


Fig. 1: The figure shows the resource utilization of a recent multi-threaded sampled simulation technique, LoopPoint, for the SPEC CPU2017 benchmarks with the `ref` inputs running eight OpenMP threads. The graph on the left shows the time required to generate the profiling data (with checkpoints stored as pinballs [56]), whereas the graph on the right shows the amount of storage required.

performance characterization. Lau et al. [47] showed a strong correlation between BBVs and region performance. Apart from BBVs, Shen et al. [63] introduced LRU stack distance vectors (LDVs) [50] to summarize program behavior for different regions. BarrierPoint [17] combines BBVs and LDVs into a signature vector (SV) in an attempt to represent more accurate features of multi-threaded applications. Furthermore, Cotson [6] and Dynamic sampling [36] record statistics such as the number of instructions executed, memory accesses, exceptions, bytes read or written, etc., in order to plot the feature of a given region. Unfortunately, none of these offline techniques are capable of representing runtime optimizations adopted by applications.

Overheads. Figure 1 illustrates the overhead of profiling data for LoopPoint (the evaluation was performed using the LoopPoint tools [49]) methodology, indicating that profile-driven methodologies incur significant overheads. When it is required to emulate an architecture (for example, simulating ARM or RISC-V binaries on x86) during profiling, it is necessary to resort to functional simulation to gather feature vectors, which can be a time-consuming process. For instance, Sandberg et al. [62] demonstrated that it took up to a month to generate profile data for SPEC CPU2006 benchmarks using simulators like `gem5`. In addition, profiling for asymmetric hardware, such as the *big.LITTLE* cores is challenging as the operating frequency (and other dynamic hardware settings) of each core is unknown at the time of profiling. Handling and storing simulation checkpoints can be a daunting task. For example, x86 checkpoints like ELFies [57] require a significant amount of storage space. These checkpoints are usually specific to the simulator being used, and they are often tied to a particular software version or hardware configuration.

Hardware and Software Dynamism. Researchers have introduced several dynamic optimization techniques in hardware and software to achieve higher performance and reduce power consumption. Techniques such as dynamic voltage and frequency scaling (DVFS) and cache reconfiguration have

been developed to adjust the hardware state in response to executed instructions and active processes. Software optimization techniques [22], [46], [67], [70] generate high-performance instructions at runtime. Additionally, dynamic scheduling techniques [29] have been developed for multi-threaded applications. In such cases, profile-driven sampling methodologies result in different performances for each execution. Methodologies such as trace-based simulations [14] or deterministic replay platforms [58] can guarantee consistent performance across multiple executions but demand extensive profiling and large storage resources. Dynamic hardware events, such as changes in core frequency, cache size, etc., are unknown during profiling. These events are performance or power-dependent and are usually hard to predict. Sherwood et al. [65] utilize a Markov Predictor to predict the phase behavior at runtime. Kihm et al. [44] propose switching to the detailed simulation mode whenever the BBV variance exceeds a specified threshold. But these methods work only for single-threaded applications as the phase behavior of synchronizing multi-threaded applications varies frequently due to the interaction of threads.

Requirements for Fast and Accurate Simulation. Sampled simulation without upfront analysis is promising under these dynamic software and hardware constraints. Therefore, it is imperative to leverage the best aspects of SimPoint-like and SMARTS-like methodologies to achieve optimal simulation efficiency and accuracy. In this work, we incorporate application analysis to guide sampled simulations, similar to SimPoint-like methodologies but without the need for upfront pre-processing, as seen in SMARTS-like methodologies. Instead, we make intelligent simulation decisions through online learning. Moreover, the proposed methodology can accommodate hardware state changes, software features, and other factors that affect simulation results. It is essential to implement efficient and lightweight online profiling, clustering, and warmup techniques for optimal performance. Therefore, to quickly estimate the performance of multithreaded applications running on next-generation dynamic hardware and software, a sampled simulation methodology is needed that can dynamically adapt to changes in the system at runtime while accurately determining relevant performance metrics.

III. THE PAC-SIM METHODOLOGY

In this section, we describe our proposal for an end-to-end sampled simulation methodology, Pac-Sim (depicted in Figure 2), that supports both dynamic hardware and software without requiring up-front workload analysis. Pac-Sim consists of five main stages: Marker Detection, Region Profiling, Clustering, Prediction, and Simulation, which are all carried out online. We have carefully designed each of these stages to minimize the runtime overhead of the methodology while maintaining the sampling accuracy. An important advantage of an online sampled simulation methodology like Pac-Sim is its ability to accurately determine the execution profile of an application without relying on the reproducibility of a program’s execution paths. This characteristic allows Pac-Sim

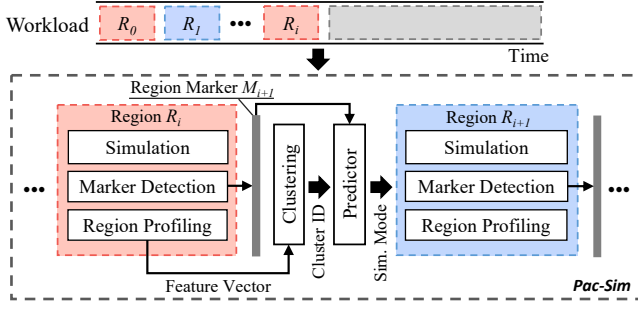


Fig. 2: Figure depicts the workflow of Pac-Sim. Consider a multi-threaded workload with regions till R_i are identified (as shown above). First, Pac-Sim monitors the application code structure to determine an appropriate region marker M_{i+1} , which marks both the end of the region R_i and the start of the region R_{i+1} . Next, the feature vector and simulation results for R_i are collected, and a prediction mechanism determines the simulation mode for region R_{i+1} . Finally, region R_{i+1} will be simulated, either in detail or in fast-forward mode.

to accurately analyze and evaluate dynamic multi-threaded applications, accounting for any performance variability that may occur at runtime.

Pac-Sim operates by making use of the program structure and runtime hardware state to identify the regions and their boundaries online. Each of these region boundaries or *markers* defines the ending of the current region and the beginning of the next region (Section III-A). Once a marker is identified, Pac-Sim collects the profiling data and simulation results of the current region (Section III-B) and clusters it with the previously identified regions to determine its cluster ID (Section III-C). This cluster ID is added to the program execution history, which is then used by the *Predictor* (Section III-D) along with the current marker and hardware state to predict whether the next region needs to be simulated in detailed mode or fast-forward mode.

While we only demonstrate the effectiveness of Pac-Sim in estimating the performance of synchronizing multi-threaded workloads in this work, our methodology has the potential to support a variety of modern workload classes, such as cloud and mobile applications, and could also be implemented for full system simulations. However, in such cases, various factors must be taken into consideration, such as kernel and driver performance, which can significantly impact the overall efficiency of the workloads. In this work, we focus on user-space workloads, and enabling support for the above-mentioned use cases is out-of-scope in this context which we leave for future work.

A. Online Region Detection

Previous research [17], [48], [60] has shown that certain program constructs, such as barriers or loops, can be utilized to characterize the phase behavior of multi-threaded applications by splitting them into a series of individually analyzable regions. Since barriers represent the global synchronization

points within a program execution, all threads align at these points, making them natural boundaries for application regions. However, relying solely on barriers to split an application may not be ideal, especially in the presence of large inter-barrier regions, as this can lead to low simulation speedups as representatives can still be too large to complete detailed simulation in a reasonable amount of time. In contrast, loops offer a finer level of granularity, allowing for greater control over the size of regions. Typically, multi-threaded applications consist of both loops and barriers in varying proportions. The online *Marker Detector* combines both of these program constructs to effectively split multi-threaded applications into regions with sizes that are well-suited for clustering while also avoiding aliasing [18]. The Marker Detector uses the following approach in order to identify the barrier- and loop-based markers online:

Barriers. Typically, a multi-threaded region begins with a `fork` call, which spawns additional worker threads and ends with a `join` call, which terminates the current thread and synchronizes with other threads. A new region is triggered at events of thread creation and termination, as regions with different active threads have different performances. For multi-threaded programs that use the OpenMP library, special function names are generated depending on the compiler used. We utilize this information in the online Marker Detector to quickly and efficiently detect barriers without much overhead.

Loops. Both loop and conditional statements use conditional branch instructions, with the target address usually given as an offset from the instruction pointer. The key difference between the two statements is that the offset of the branch instructions in a loop statement is usually negative, whereas that in a conditional statement is positive. While there are exceptions, it is generally sufficient to select conditional branches with negative offsets as markers for loops. We also make sure to disregard spinloops from our analysis.

As an application executes, the Marker Detector identifies markers online, splitting the application into multiple regions. While doing so, it also monitors the region sizes to ensure they fall approximately within the bounds of δ_{min} and δ_{max} instructions. A minimum number of instructions, δ_{min} , is necessary to capture the frequent variations in the multi-threaded program behavior and accurately cluster the obtained regions. Whenever the Marker Detector chooses barrier-based markers as region boundaries, the size of the region can be as small as δ_{min} instructions but no larger than δ_{max} instructions. Otherwise, the Marker Detector chooses the first loop-based marker it encounters beyond δ_{max} instructions as the next region boundary. For loop-bounded regions, it is necessary to keep region sizes large enough to avoid aliasing [16]. In our experiments with fixed region sizes of 10 million, 20 million, 50 million, and 100 million instructions, the SPEC CPU2017 benchmarks showed average error rates of 6.9%, 3.3%, 1.8%, and 1.8%, respectively. We set the lower bound δ_{min} to be 20 million to ensure sampling accuracy and the upper bound δ_{max} to be 50 million for better performance.

Hardware State. The Marker Detector also monitors the

hardware state of the simulated system. If it detects changes, the current region is ended at the next marker so that each region has a consistent hardware state. Once a marker is detected, the program counter (PC) and the hardware state of the simulated system are collected and stored corresponding to the marker. The collected hardware state includes the system parameters, like processor frequency, cache size/configuration, power management techniques, etc., that can be configured during runtime.

B. Online Region Profiling

Conventionally, BBVs have been used to characterize the execution behavior of code regions, as they have been shown to exhibit a strong correlation with the region's performance [47]. BBVs record the execution counts of each basic block (i.e., code blocks with single entry and exit points) within a given code region. The number of dimensions for a BBV depends on the number of basic blocks executed, which could range anywhere from thousands to even millions for very large applications. This presents a major challenge for online analysis of BBVs as the time and effort required for this stage would significantly increase as the vector dimensionality increases. SimPoint [64] uses random linear projection [23] to overcome this problem. However, this method is not suitable for our online algorithm as the matrix-vector multiplication operations involved could introduce significant runtime overheads.

To overcome these issues, we propose a fast online BBV generation technique (illustrated in Figure 3). Rather than creating a fixed-size BBV for each region, we use an online projection technique to generate fixed-size vectors BBV'_i for each basic block BB_i , where the elements of BBV'_i are computed by multiplying the instruction count of a basic block with the hash results of its program counter (PC) value. We use the hash function `drand48()`, which generates pseudo-random numbers for an integer value input. The initial four dimensions of the online BBV are determined using the hash values utilizing inputs PC, PC+1, PC+2, and PC+3, respectively. The values of the subsequent four dimensions are generated using the output of the preceding four dimensions as inputs to the hash function. We experimentally determined that using 16 dimensions adequately captures the representation of a region using the online BBV. The resultant BBV'_i vectors are then accumulated to obtain the per-thread BBV (BBV'_{online}) for the given region, which can be represented as:

$$BBV'_{online} = \sum_i BBV'_i = \sum_i (BBV_i \cdot M_{proj}),$$

where the values of the elements in M_{proj} are generated using hash functions as mentioned above. This BBV'_{online} for a region is analogous to the BBV utilized in SimPoint, which is obtained through random linear projection. The projected down BBV used in SimPoint, $BBV'_{offline}$, is obtained from the dot product of the actual BBV of the region and projection matrix M_{proj} :

$$BBV'_{offline} = BBV \cdot M_{proj} = \sum_i (BBV_i \cdot M_{proj}).$$

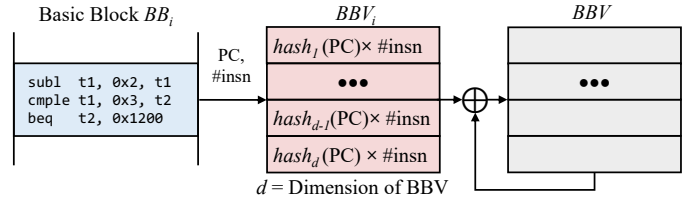


Fig. 3: The figure shows the workflow of online BBV generation. Whenever a basic block BB_i is encountered, a corresponding execution fingerprint BBV'_i is generated using hash functions applied to the program counter of BB_i and the number of instructions it contains. $hash_1$ to $hash_d$ are d distinct hash functions, where d is the dimension of the BBV. The BBV for each region is obtained by accumulating all BBV'_i s that belong to the region.

We then normalize these per-thread BBVs and concatenate them into a single global-BBV vector to represent the software feature of a given multi-threaded code region. This eliminates the need to perform computationally intensive dimensionality reduction techniques online and allows Pac-Sim to quickly determine the BBVs without much overhead.

C. Determining Region Similarity

Pac-Sim employs an online clustering mechanism to group regions with similar execution behavior based on the feature vectors collected for each region in the online profiling stage. The clustering, which is done at the end of simulating each region, is required for the learning process of the Predictor. Prior works, like SimPoint [64], cluster feature vectors using the k-means algorithm [37]. However, k-means uses an iterative refinement technique that is computationally intensive, and therefore, it would not be practical to use this algorithm for determining region similarity online.

In order to reduce this computational load and enable real-time clustering, we devise an alternative technique for clustering feature vectors (i.e., global-BBVs) in Pac-Sim. In our technique, we maintain two separate queues: (i) detailed queue and (ii) fast-forward queue. The detailed queue includes the BBVs corresponding to the regions that have been simulated in detail, while the fast-forward queue includes those corresponding to the regions that have been fast-forwarded. When a new BBV is recorded, it is first compared with the BBVs in the detailed queue. If its distance from any of these BBVs is less than the specified threshold θ , then we return the cluster ID of the closest region. If there is no region whose distance is less than θ , we repeat the same procedure with the regions in fast-forward queue. If we still don't find similar regions, we assign a new cluster ID for the current region and insert it into the BBV queue corresponding to its simulation mode. In our experiments, we set $\theta = 0.05$ to ensure a reasonable simulation accuracy while maintaining high speedups. To further improve the efficiency of our clustering technique, we incorporate the triangle inequality optimization [32] into our algorithm, which

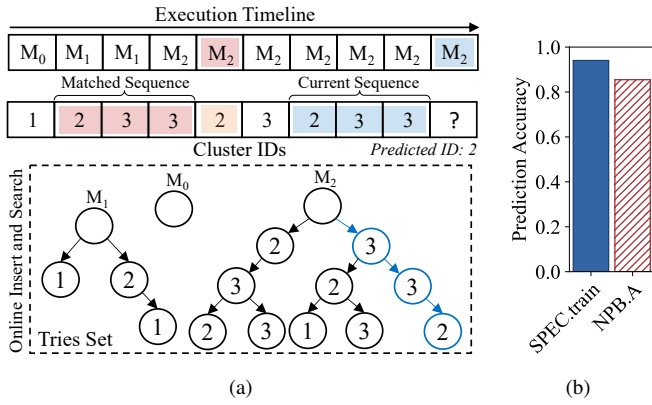


Fig. 4: The predictor utilizes the trie [12] data structure to quickly predict the cluster ID of the next region by searching for a similar history with the same region start marker M_i . In this example, the cluster ID of the next region is predicted to be 2 since the prior region with the cluster ID of 2 has the same start marker M_2 and the longest matching sequence ($2 \rightarrow 3 \rightarrow 3$). Plot (b) shows the accuracy of the predictor for different benchmark suites.

can skip redundant BBV distance calculations. We consider Euclidean distance for all BBV distance calculations.

D. Prediction Mechanism

Pac-Sim employs a *Predictor* – an online prediction mechanism that leverages region markers, execution history, and hardware state to predict the phase behavior of the next region in an application and decide its simulation mode at runtime.

Region Markers. The Marker Detector identifies PC-based region markers that act as the boundaries of the regions. In certain cases, using region markers to classify regions is effective for applications where the same part of the code displays similar phase behavior, as in the case of `619.lbm_s.1` and `644.nab_s.1` using train inputs.

Execution History. When executing the same part of the source code, differences in memory access patterns, branching, etc., can result in varying phase behavior at runtime. We, therefore, make use of execution history, which is a sequence of the cluster IDs of prior regions, to predict these differences in the phase behavior among applications.

Hardware State. Pac-Sim takes into account the state of the simulated system, such as core frequency, while executing each region. *Predictor* predicts the next region to be detailed mode if there are no prior similar regions with the same hardware state. The Predictor decides the cluster ID of the next region by choosing the cluster ID of the previous region with the same region marker and has the longest matching sequence. For the regions that do not have a previous region with the same start marker or the same history, Pac-Sim enables detailed execution for that region. Then it decides the simulation mode of the next region by checking whether prior regions with cluster ID and the current hardware state are

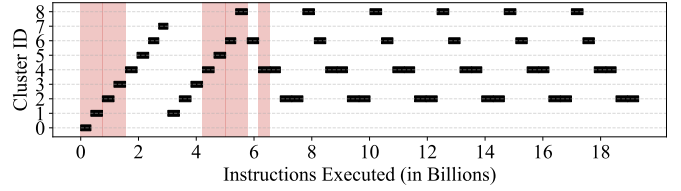


Fig. 5: The graph shows the regions identified using Pac-Sim for the NPB benchmark `ft`, grouped together with the respective cluster they belong to. The shaded portion represents the regions that are simulated in detail.

simulated in detailed mode. This history is learned online and is updated every time Pac-Sim finishes simulating a region.

To accelerate this stage for large application lengths, we further optimize our clustering algorithm to reduce its average search time complexity from $\mathcal{O}(n^2)$ to $\mathcal{O}(n)$. This is achieved by maintaining the execution history in a *trie* [12] data structure, with a maximum depth of 16, which allows for more efficient *search* and *insert* operations. In Pac-Sim, we utilize the trie data structure to maintain the execution history of the application being simulated and quickly predict the cluster ID of the next region based on this information.

Figure 4a illustrates the usage of tries to predict the cluster ID of the next region by considering the example of a hypothetical execution sequence. *Insert*: The cluster ID of the current region is inserted into the trie. In Figure 4a, when the online clustering of the fifth region is finished, we insert the current cluster ID 2 for both branches corresponding to the three histories: $3, 3 \rightarrow 3$, and $2 \rightarrow 3 \rightarrow 3$. *Search*: Once Insert of the current region is completed, the cluster ID of the next region is predicted by searching the trie for a matching cluster ID sequence. The search operation ends when the sequence matches one of the leaf node paths. Note that two regions having the same marker do not necessarily mean that the regions belong to the same cluster.

Figure 4b shows the average accuracies of the online predictor for the benchmarks of SPEC CPU2017 and NPB are 94% and 85%, respectively, ensuring the sampling accuracy and performance of Pac-Sim. The accuracy of the predictor is determined by comparing the predicted cluster ID prior to simulating the region with the actual cluster ID obtained through clustering after simulation. Figure 5 shows the results of the Predictor in clustering different regions identified by Pac-Sim simulating the `ft` benchmark from the NPB benchmark suite using eight threads. We observe that the majority of regions from each cluster are simulated in detail (shaded portions). This is in accordance with the learning phase of our algorithm where Pac-Sim works to establish a comprehensive understanding of the phase behavior of the application.

E. Simulation by Application Reconstruction

Previously proposed multi-threaded sampling methodologies [17], [60] rely fully on offline analysis to determine the regions that need to be simulated in detail. Pac-Sim assumes

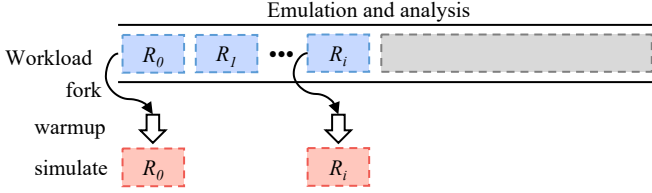


Fig. 6: The workflow of Pac-Sim when the representative regions are simulated in parallel. Pac-Sim starts in the emulation mode, collecting feature vectors and MTR [8] warmup data online, and then predicts the simulation mode of the next region. For regions predicted for detailed mode, Pac-Sim forks new processes to perform warmup and detailed simulation.

no prior knowledge about the nature of the workload that it is about to simulate. Instead, it (a) samples regions online during the simulation and (b) uses the detailed simulation results of previous regions to estimate the performance of the current fast-forwarded region by applying the four different methods described below successively until convergence is reached.

- i. Use the detailed performance metrics of a region that belongs to the same cluster and has the same start marker as the current region.
- ii. Use the detailed performance metrics of a region that belongs to the same cluster.
- iii. Use the performance of a region with the closest BBV ($\theta = 1$) and the same number of active threads.
- iv. Use the average performance of the regions that have the same number of active threads.

We extrapolate the runtime T_r of a fast-forwarded region r using the region r' by $T_r = T_{r'} \frac{insn_r}{insn_{r'}}$, where $T_{r'}$ is the runtime of the previous region r' identified above, and $insn_r$ and $insn_{r'}$ are the maximum instruction counts among all threads for the regions r and r' , respectively.

Runtime Hardware Events. Pac-Sim takes into account the state of the simulated system while estimating the performance of the fast-forwarded region. As runtime hardware events can happen at any time, we do not guarantee the regions that are divided by those events to be large enough. In such cases, we estimate the performance of these regions using the closest previous region with the same hardware state, as these regions are too small to be clustered. Moreover, the impact of these regions on the overall application performance is typically negligible as the regions are too short.

F. Sampled Simulation in Parallel

Pac-Sim is primarily targeted for runtime varying scenarios using live sampling. However, for statically scheduled multi-threaded applications, Pac-Sim can support sampled simulation in parallel, similar to checkpoint-based mechanisms, to further speed up the sampled simulation. The workflow of Pac-Sim for parallel simulation is shown in Figure 6. Previous methods, like LiveSim [41] and LoopPoint [60], require offline analysis and store checkpoints for sampled simulation. A huge amount of storage is required for these methods, as mentioned

in Section II. Pac-Sim starts in emulation mode, collecting feature vectors and warmup data online, and then predicts the simulation mode of the next region. For regions predicted for detailed mode, Pac-Sim forks new processes, which run in parallel, to perform warmup and detailed simulation. Pac-Sim reconstructs the performance of the entire application once the whole application is emulated and the simulation of all regions is completed.

G. Microarchitectural Warmup

One of the major challenges of sampled simulation is to build up the accurate microarchitectural state prior to the detailed simulation of each region. Choosing the right warmup technique that can directly build this state is crucial in order to achieve the highest speedup. Methodologies like SMARTS [69] and time-based sampling techniques [5], [16] keep functional warming enabled for the entire sampled simulation leading to large slowdowns. We find that the statistical warmup techniques [8], [40], [54], [66] can reconstruct the accurate microarchitectural state of a simulated system online. We select MTR [8] to be used with Pac-Sim as it can rapidly collect memory reference patterns during fast-forward mode and reconstruct the cache state before switching to detailed mode. Caches are larger structures compared to branch predictors or prefetchers, and hence we limit the simulation infrastructure to explicit cache warming, as the smaller structures, like prefetchers, are warmed quickly. Moreover, we maintain larger regions (minimum region sizes of 20 million instructions) to achieve good warm-up performance for other microarchitectural structures. It is also possible to increase the amount of warmup needed for different structures, and there are different ways to solve this problem, but warmup is a challenge that is workload-specific. The exploration of additional warmup scenarios is outside the scope of this work.

IV. EXPERIMENTAL SETUP

In this section, we describe the experimental setup used to evaluate Pac-Sim. We begin by providing the details of the simulation framework used in our experiments. We then describe the different workloads that are used to evaluate the performance of our methodology.

A. Simulation Tools

In this work, we use a modified version of the Sniper multi-core simulator [15] (version 7.4), which is updated to support loop-based and barrier-based region specifications in order to evaluate Pac-Sim. Sniper is a many-core simulator using high-level abstract models and is widely used for architectural evaluation and design space exploration. Note that our methodology does not utilize any features specific to the Sniper simulator. Therefore, porting the methodology to other simulators, such as gem5 [10] or ZSim [61], should be relatively straightforward. To demonstrate that Pac-Sim is indeed a microarchitecture-independent methodology, we experimentally evaluate it by running simulations upon two different processor configurations that mimic the performance/behavior

of Intel’s Gainestown and Skylake² [28] microarchitectures using Sniper. The configuration details for each of these models are listed in Table I.

TABLE I: The configuration parameters we used for Gainestown and Skylake microarchitectures on Sniper.

Component	Gainestown Parameters	Skylake Parameters
Processor	1, 8 cores	1, 8 cores
Core	2.66 GHz, 128-entry ROB	2.66 / 3.7 GHz, 224-entry ROB
L1-I / L1-D	32KB, 4 / 8 way, LRU	32KB, 8 / 8 way, LRU
L2 cache	256KB, 8 way, LRU	1MB, 16 way, LRU
L3 cache	8MB (shared), 16 way, LRU	22MB (shared), 12 way, LRU

In order to speed up the simulation, Pac-Sim intelligently switches among the three simulation modes supported by Sniper, namely, fast-forward mode, cache-only mode, and detailed simulation mode. The fast-forward mode is used to reach a particular point in an application during simulation without enabling the performance models. The cache-only mode performs the functional warming of the caches, whereas the detailed simulation mode is the default simulation mode that enables the timing model for performance estimation. For Pac-Sim, we capitalize on this split execution and timing model architecture to fast-forward in the front-end of the simulator so that the simulation wall time is further minimized.

Every time Pac-Sim switches from fast-forward to detailed simulation mode; the cache state is reconstructed at the beginning of the region using the memory time-stamp record (MTR) [8] technique. We implement MTR in Sniper to collect the cache line information accessed by each *Load* and *Store* instruction during simulation, ordered in LRU fashion per set, and then inject the requests into the cache in the correct order to rebuild the appropriate cache state.

B. Benchmarks Used

To demonstrate the wide applicability of Pac-Sim, we experimentally evaluate the methodology using multiple benchmark suites such as (i) the SPEC CPU2017 benchmark suite [13], (ii) the NAS Parallel Benchmarks (NPB) [7] version 3.4.2, and (iii) the PARSEC [9] version 3.0 benchmark suite. Note that these are multi-threaded benchmarks that synchronize frequently and share memory.

We configure these benchmarks to use two different multi-threaded programming models, namely OpenMP [55] and OmpSs [29]. OpenMP [55] provides a set of compiler directives, library routines, and environment variables that help developers to parallelize their code. On the other hand, OmpSs [29] extends OpenMP, and it is able to dynamically manage and schedule tasks to maximize multi-threaded application performance. We set up the multi-threaded benchmarks to use *passive* thread wait policy, meaning that the threads will sleep while waiting for other threads at a synchronization point.

²Note that Gainestown is the latest microarchitecture available on Sniper simulator that has been validated against hardware. We made modifications to the back-end of Sniper to support the contention model and instruction latencies for Skylake architecture.

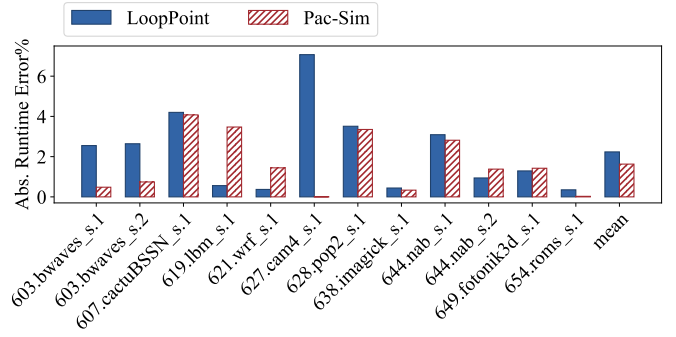


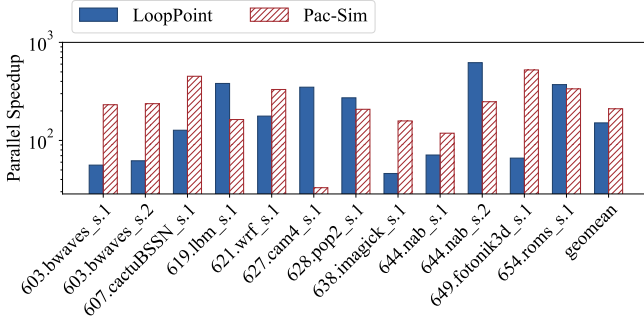
Fig. 7: A comparison of the absolute runtime prediction error for LoopPoint and Pac-Sim for 8-threaded SPEC CPU2017 benchmarks using train inputs. Pac-Sim achieves similar levels of accuracy compared to LoopPoint.

SPEC CPU2017 is a collection of benchmarks used for performance evaluation in computer architecture research. In our experiments, we use the *speed* version of multi-threaded SPEC CPU2017 benchmarks that are parallelized with OpenMP. The benchmarks are compiled using GCC 6.4.0 and GFortran with the `-O3` compiler flag for x86-64 architecture. We configure these benchmarks to run with eight threads and evaluate them using the *train* input set. NAS Parallel Benchmarks (NPB) [7] is another set of benchmarks widely used to evaluate the performance of highly parallel systems in computer architecture. The reference implementations of these benchmarks are available in the two most commonly used programming models, i.e., MPI and OpenMP. In our experiments, we use the OpenMP-based implementation with input *class A* and generate the binaries using `icc` compiler (with `-O2` flag) as part of the Intel oneAPI (version 2022.0.2) toolkit. We also present experimental evaluations of Pac-Sim using PARSEC, which is another standard benchmark suite consisting of computationally intensive applications designed to facilitate the study of multi-core systems with shared memory. PARSEC implementations are available in both OpenMP and OmpSs [20] versions. In our experiments, we use both these versions with the *simlarge* input set.

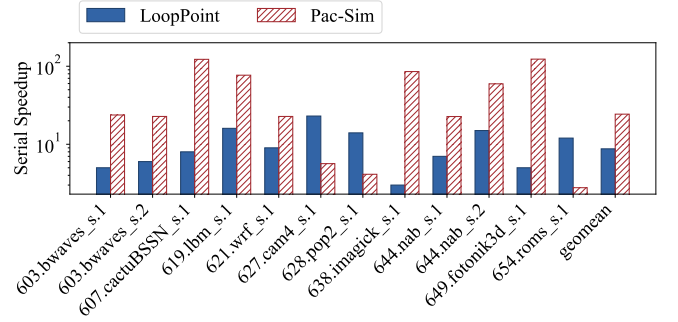
V. EVALUATION

In this section, we first present a comprehensive evaluation of Pac-Sim, comparing its efficacy with the current state-of-the-art. Additionally, we provide experimental evidence showing that Pac-Sim is indeed a hardware-independent methodology. Finally, we present case studies that demonstrate the applicability and effectiveness of Pac-Sim in estimating workload performance in dynamic, multi-threaded hardware and software environments. Note that, throughout this paper, the term runtime refers to the simulated runtime of the application, whereas the term wall-time refers to the actual time taken by the simulator to finish the run.

Evaluation metrics. In order to evaluate the effectiveness of any simulation methodology, it is crucial to quantitatively measure its performance in terms of two critical metrics:

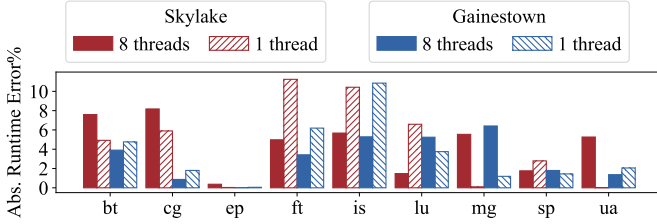


(a) Parallel Speedup

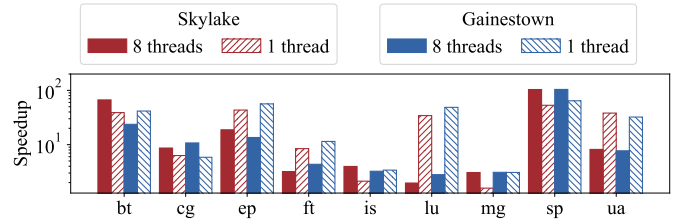


(b) Serial Speedup

Fig. 8: The parallel and serial speedups of Pac-Sim are compared with that of LoopPoint for 8-threaded SPEC CPU2017 benchmarks using train inputs. For speedup calculations, the simulation walltime corresponding to Pac-Sim includes both online analysis and simulation time, whereas, for LoopPoint, we consider only the checkpoint simulation time, excluding the time required for offline profiling and checkpoint generation.



(a) Accuracy



(b) Serial Speedup

Fig. 9: The accuracy and serial speedup achieved for Pac-Sim methodology when simulated using Gainestown and Skylake architectures for NPB benchmarks with class A inputs running eight threads and one thread.

accuracy and *speedup*. In our experiments, we define these metrics in the following manner:

Accuracy: We assess the accuracy of our proposed methodology by comparing the simulation runtime obtained from the full simulation and the sampled simulation in terms of absolute runtime prediction error Δ_{time} , which is defined as

$$\Delta_{time} = \frac{|T_{full} - T_{sample}|}{T_{full}},$$

where T_{full} represents the simulation runtime obtained from the full run, and T_{sample} represents the simulation runtime extrapolated from the sampled simulation. It is important to note that in our evaluation, we use the runtime (execution time as inferred from simulation) of the application as the performance metric to measure the accuracy of sampling. This is because time-per-program is the gold-standard performance measure, and IPC is not a valid performance metric for multi-threaded applications [3].

Speedup: In our experiments, we calculate the speedup by taking the ratio of the wall-clock time for the full simulation to that of the sampled simulation and the average speedup by computing the geometric mean of the speedups across all benchmarks. Serial speedup is defined as the speedup achieved when all representative regions are simulated sequentially,

while parallel speedup is obtained when the representative regions are simulated in parallel, assuming infinite resources.

A. Comparison with the State-of-the-Art

In this section, we evaluate the performance of Pac-Sim in comparison to the state-of-the-art profile-driven sampled simulation methodology, LoopPoint [60]. While several other profile-driven methodologies exist, LoopPoint provides the benefit of being applicable across a variety of application and synchronization types. It has also been shown to outperform other multithreaded sampled simulation methodologies (such as BarrierPoint) in terms of speedup and accuracy, thus serving as a strong baseline for our evaluations. We now report the results of our simulation experiments evaluating and comparing the performance of these two methodologies using the SPEC CPU2017 benchmarks.

Accuracy. Figure 7 shows a comparison of absolute runtime prediction errors for Pac-Sim and LoopPoint obtained for the 8-threaded SPEC CPU2017 benchmarks using train inputs. Our analysis reveals that, in most cases, Pac-Sim performs comparably with LoopPoint in predicting the runtime of the applications, with the individual errors differing by no more than 2 to 3%. The relatively higher errors for some applications, such as 619.lbm_s.1 is because Pac-Sim relies

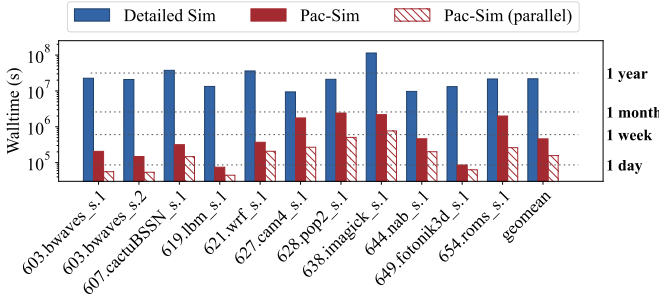


Fig. 10: A comparison of the estimated walltime for fully detailed simulation and sampled simulation using the serial and parallel versions of Pac-Sim for 8-threaded SPEC CPU2017 benchmarks using ref inputs. The estimated walltime includes the time required for online analysis, warmup, and simulation.

on online extrapolation to estimate application performance using the limited profile data that is available from regions that have already been simulated. Whereas methodologies like LoopPoint rely on offline profiling and therefore utilize the information about the whole application.

Speedup. Figure 8 shows the speedup comparison of Pac-Sim and LoopPoint for the SPEC CPU2017 benchmarks using train inputs running eight threads. Figure 8a shows the parallel speedup for which Pac-Sim outperforms LoopPoint in most cases (7 out of 12 benchmarks). The primary reason for this is that Pac-Sim uses smaller regions as compared to LoopPoint. Although Pac-Sim requires emulation of the entire application, the online analysis overhead is minimized, and therefore, the average parallel speedup for SPEC CPU2017 benchmarks (train inputs) using Pac-Sim is $210.3\times$, which is larger than that obtained for LoopPoint ($150.97\times$).

Figure 8b shows the serial speedup, and we observe Pac-Sim outperforms LoopPoint in most cases, attaining a maximum serial speedup of $123.32\times$. While the online analysis can introduce some runtime overheads, the performance advantages of Pac-Sim seem to outweigh these overheads in most cases. However, there are some cases where LoopPoint performs better than Pac-Sim, such as for `627.cam4_s.1` and `628.pop2_s.1` benchmarks in Figure 8b. This is mainly because Pac-Sim uses a small clustering threshold (0.05) for the online clustering in order to maintain high accuracy.

Efficacy in Evaluating Realistic Workloads. The full detailed simulation of SPEC CPU2017 benchmarks with reference inputs takes an extremely long time – about a year on average using multi-core simulators like Sniper. Instead, we estimate their simulation walltime by considering the instruction count of the benchmark using reference inputs along with the average simulation rate of the benchmark using train inputs. The walltime of Pac-Sim includes the time required for online analysis and emulation of the entire workload along with the time for detailed simulation of the representative regions. Figure 10 shows that Pac-Sim takes less than a week, on average, to run the entire application sequentially, while the parallel version of Pac-Sim takes about 1.8 days on average. In

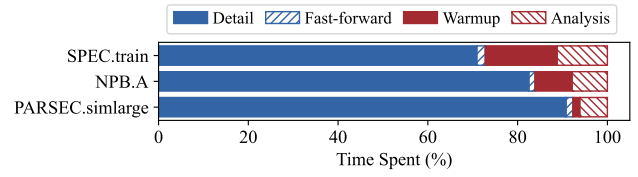


Fig. 11: The graph shows the percentage of time that Pac-Sim spends at each phase during the sampled simulation of each benchmark suite (average across all benchmarks). The Analysis part includes online marker detection, region profiling, clustering, and prediction.

experiments where the microarchitecture structures like cache size are adjusted or when the application itself undergoes instruction-level modifications, it is necessary to regenerate the checkpoints. In such cases, Pac-Sim is more appropriate as LoopPoint takes 6.2 days on average (shown in Figure 1) to complete its preprocessing before simulation.

Microarchitecture-agnostic sampling. In addition to achieving high accuracy and speedups, Pac-Sim also provides the advantage of being a microarchitecture-independent methodology. We experimentally demonstrate this by evaluating our methodology with two different processor configurations, namely the Gainestown and Skylake microarchitectures, for the NPB benchmarks that run using one thread and eight threads. The accuracy and speedup numbers obtained in our experiments are plotted in Figure 9a and Figure 9b, respectively. From Figure 9a, we can observe that the absolute runtime errors estimated by Pac-Sim for all NPB applications are quite low (all under 8%) and are similar for both these processor configurations (differing by 5% at most). Moreover, the speedups obtained for both configurations are similar for most benchmarks, as observed in Figure 9b. Hence, the choice of a target microarchitecture for evaluation does not affect the efficacy of Pac-Sim.

Wall-time Distribution. We show the time spent by Pac-Sim in different stages of sampled simulation. Figure 11 shows the average time spent in the online analysis stage for NPB (class A inputs) and SPEC CPU2017 (train inputs) benchmarks is 7.88% and 11.20%, respectively. This is the result of the optimizations described in Section III, which are applied to the analysis part. Moreover, Pac-Sim spends 8.25% and 16.00% of the execution time on warmup for NPB and SPEC CPU2017 benchmarks, respectively. This is because Pac-Sim needs to reconstruct the memory access patterns at the beginning of the detailed simulation of a region. Note that Pac-Sim brings down the time spent in profiling/analysis of the benchmarks significantly as compared to prior profile-driven methodologies for sampled simulation.

B. Case Studies

We showcase the versatility of Pac-Sim through several compelling case studies. Firstly, we demonstrate that our methodology remains agnostic to dynamic thread scheduling decisions made during runtime, highlighting its robustness

and adaptability. Next, we provide examples of how Pac-Sim operates seamlessly in the presence of various runtime hardware events, further cementing its reliability. Finally, we exhibit the applicability of the proposed methodology in hardware-software co-design studies, showcasing its potential to facilitate more efficient and effective design processes.

1) Dynamically Scheduled Software

With the advent of multi-core and many-core processors, efficient parallel execution of dynamically scheduled multi-threaded applications has become crucial to the performance of modern computing systems. However, the non-determinism resulting from the execution of such applications on multi-core platforms often leads to notable performance variability across multiple runs. At the software level, such variability could arise from dynamic scheduling of system jobs, thread migrations, load balancing optimizations, or contention on shared resources at runtime.

TABLE II: Table shows the IPC of `freqmine` benchmark from the PARSEC benchmark suite using *simlarge* input for threads 0, ..., 7. Pac-Sim shows the details of dynamically scheduled software whose IPC and thread mapping differ across two runs.

Thread ID	0	1	2	3	4	5	6	7	Aggr.
IPC _{run1}	0.15	0.09	1.75	0.43	0.07	0.07	0.10	0.09	2.75
IPC _{run2}	0.15	0.09	1.76	0.07	0.44	0.07	0.09	0.10	2.76

Table II illustrates the thread-level differences in terms of IPC for two runs of the OpenMP-parallelized `freqmine` application from the PARSEC benchmark suite. There are variations in the per-thread IPCs between the two runs, particularly for thread IDs 3 and 4. To investigate the impact of these variations on conventionally used sampling techniques, we conducted two profiling runs of `freqmine` using LoopPoint. The experimental result revealed that 14% of the regions were clustered differently for the second run as compared to the first. This presents a challenge for sampled simulation, which relies on profiling data from a prior execution to guide simulation in subsequent runs, as dynamic applications have variant profiling data across different executions. To overcome these issues, Pac-Sim profiles and clusters the regions online and simulates them in the same run, thereby accounting for any performance variability that may occur at runtime.

To demonstrate the effectiveness of Pac-Sim in this regard, we now present an experimental study of dynamically scheduled multi-threaded versions of PARSEC with *simlarge* inputs and NPB with class A inputs. While the per-thread behavior varies for dynamically scheduled applications, the global execution time and global IPC remain consistent across multiple runs. In Table II, we observe that while there are some variations in per-thread behavior, the aggregate IPCs across the two runs remain nearly unchanged. Figure 12 demonstrates the average runtime prediction errors of Pac-Sim simulating dynamically scheduled multi-threaded applications. We run the benchmarks multiple times in full detailed mode and using Pac-Sim. The errors are calculated by comparing

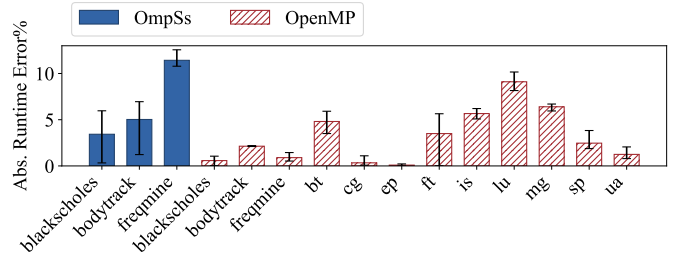


Fig. 12: Figure shows the average error rates (from five different runs) and error bars in predicting the runtime of dynamically scheduled benchmarks. We use PARSEC benchmarks with the *simlarge* input using OmpSs and OpenMP, and NPB benchmarks with *class A* inputs using OpenMP runtime.

the runtime obtained using Pac-Sim with the average runtime obtained from the full detailed simulations. The results show that Pac-Sim achieves a very low error in predicting the runtime of dynamically scheduled software (3.81% on average). The benchmark, `freqmine`, which shows the largest IPC variation maintains an average error of 11.43%. Moreover, Pac-Sim demonstrates speedups of up to $43.96\times$ ($6.29\times$ on average) for all dynamically scheduled benchmarks.

2) Dynamic Hardware Events

Dynamic event-based hardware optimizations help improve performance gains and energy efficiency in modern architectures. DVFS [33], [42], [45] is one of the most widely employed dynamic hardware event-based optimization techniques. It monitors core frequencies and load variations in order to match the system power consumption with the required level of performance by triggering voltage and frequency optimizations at runtime. These optimizations may lead to a diverse range of dynamic hardware states (i.e., core frequency, power configurations) over a given run, consequently resulting in a significant degree of performance variability for a given workload across different executions.

Pac-Sim deals with this performance variability by monitoring the simulated hardware events at runtime. While prior sampled simulation methodologies only support dynamic hardware events triggered at region boundaries to ensure that the hardware state remains constant for a given region, Pac-Sim supports hardware events at any time during the application execution. Each time an event occurs, Pac-Sim triggers a new region to ensure hardware state consistency within that region. The predictor then speculates the cluster ID of the next region and checks the execution history to determine whether similar regions (i.e., regions with the same cluster ID and hardware state) were previously encountered. If a similar region has been previously simulated in detail, the region is fast-forwarded; otherwise, a detailed simulation mode is triggered.

We now present an experimental study demonstrating the effectiveness of Pac-Sim in handling the variability caused by dynamic hardware events by specifically considering the case of DVFS-optimized workloads. In our experiments, we evaluate the performance of the benchmarks by comparing

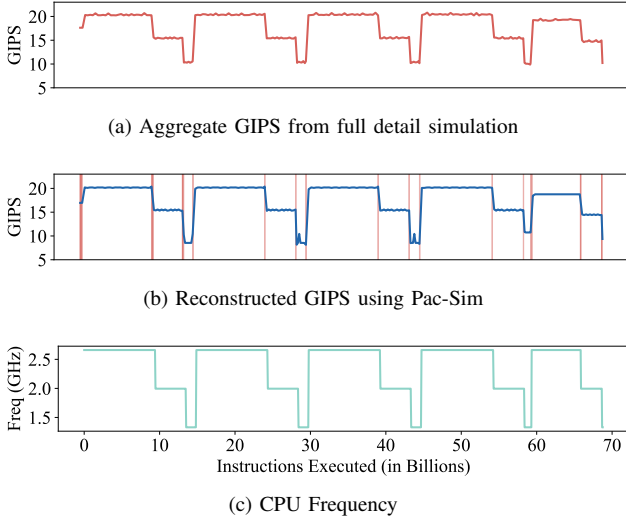


Fig. 13: The aggregate giga (billion) instructions per second (GIPS) of the full run (a), reconstructed GIPS using Pac-Sim (b), and the varying CPU frequency for all CPUs (c) 644.nab_s.1 benchmark with train inputs running 8 threads. The shaded regions in (b) represent the regions simulated in detail. The figures share the same x-axis.

the results of Pac-Sim with the baseline while changing the frequency at predetermined intervals; however, just like in actual DVFS-optimized executions, the information on the frequency changes is not available to the simulator a priori. In order to evaluate the performance of Pac-Sim, we consider a DVFS scenario in which the processor frequency f switches among a fixed range of values, i.e., $f \in \{1.33 \text{ GHz}, 2.00 \text{ GHz}, 2.66 \text{ GHz}\}$ as shown in Figure 13c.

For this scenario, we measure the aggregate giga/billion instructions per second (GIPS) values obtained from both the full detailed simulation and Pac-Sim over the entire execution. The findings of our experiment are presented in Figure 13. We observe that the GIPS values obtained from both the full simulation (Figure 13a) and Pac-Sim (Figure 13b) exhibit a great deal of similarity, indicating Pac-Sim’s effectiveness in estimating the performance of a dynamically optimized workload with a high level of accuracy. Furthermore, our findings reveal that Pac-Sim simulates only a small fraction of the entire application in detail (depicted by shaded regions in Figure 13b). Notably, most of the detailed simulation occurs either at points of change in the phase behavior of the application or hardware states. This demonstrates that Pac-Sim can use this information to identify a minimal representative subset for applications using online analysis.

3) Hardware-Software Co-design

Hardware-software co-design is an emerging field of study that optimizes the system performance by concurrently designing the compiler and hardware components of a system to exploit the synergy between the two. Prior works [35], [39], [72] have investigated several directions in this context. To identify the most effective strategies, hardware-software

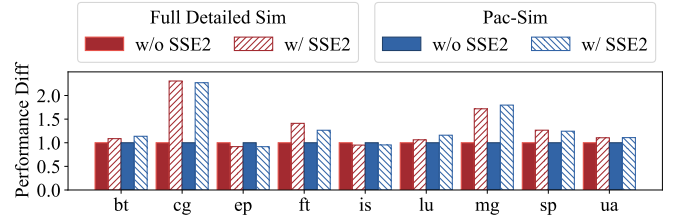


Fig. 14: The figure shows the absolute difference in performance (in terms of runtime) for NPB benchmarks using class A inputs and 8 threads with (w/) and without (w/o) SSE2 simulated in detailed mode and with Pac-Sim.

co-design research relies on fast and accurate architectural simulation methodologies to explore the design space efficiently. However, among existing methodologies, the profile-driven methodologies [17], [64] incur significant profiling and preprocessing costs, as shown in Figure 1, whereas the statistical sampling methodologies [43], [69] (which don’t rely on preprocessing) have low speedups.

Pac-Sim addresses these issues by sampling and analyzing the regions online. Thus, it incurs no additional profiling cost if new compilers are used or new applications are generated, enabling fast and efficient exploration of hardware-software co-design space. To demonstrate the effectiveness of Pac-Sim in this regard, we now present a performance evaluation study of the NPB benchmarks under different compiler optimization techniques. We study the impact of SIMD (Single Instruction, Multiple Data) optimizations on the generated binaries using both Pac-Sim and full detailed simulations. SIMD-enabled processors are equipped with special-purpose registers that can simultaneously load multiple machine words and perform operations on them in parallel in order to improve processor performance. For instance, the Streaming SIMD Extensions 2 (SSE2) instruction set uses 128-bit XMM registers to process packed data elements at once.

In our experiments, we measure the performance improvement (in terms of runtime) obtained by enabling SSE2 and compare it against the baseline. The results of our simulations are graphed in Figure 14. We observe that the average difference in the performance improvements obtained from full detailed mode and Pac-Sim is 3.65%. Specifically, Pac-Sim reveals the performance effects of SIMD instructions. For example, some benchmarks achieve a significant speedup over the baseline as these applications meet the `icc` vectorization criteria [26]. `ft` calculates a 3D fast Fourier transform, and its innermost loop consists of multiply-add statements with contiguous memory accesses and no data dependency. On the other hand, `is`, which uses the quick sort algorithm, is hard to vectorize. The SIMD overheads resulting from register transfer costs exacerbate the overall application performance.

VI. RELATED WORK

We have discussed the most relevant previous works in Section II. Sampled simulation has been an active research area

for several decades, and several techniques were proposed [5], [6], [16], [17], [38], [41], [60], [62], [64], [68], [69] in this direction for different workload classes primarily for the reduction of simulation time and resources. In order to simulate the identified sample correctly, it is important to reconstruct the microarchitectural state of the system using techniques like statistical warming [8], [30], [40], [53], checkpoint-based warming [54], [66], or by enabling functional simulation. Analytical modeling is yet another solution to evaluate a complex workload quickly. Prior works proposed analytical models to derive the performance of processors [24], [34], cache miss rates [31], branch miss rates [25], DVFS performance [1], etc. However, analytical performance modeling can be limited in supporting new designs, requiring new models for each.

VII. CONCLUSION

This work proposes a novel sampled simulation methodology and infrastructure called Pac-Sim. The work focuses on what is needed to simulate dynamic software that responds to workload- and run-time-specific execution conditions. Pac-Sim is the first, to the best of our knowledge, to propose a sampling solution that simulates these dynamic conditions in both a fast (up to $523.5\times$ speedup, $210.3\times$ on average) and accurate way (average errors of 1.63% and 3.81% for statically and dynamically scheduled benchmarks, respectively).

REFERENCES

- [1] S. Akram, J. B. Sartor, and L. Eeckhout, "DVFS performance prediction for managed multithreaded applications," in *2016 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS)*, 2016, pp. 12–23.
- [2] A. Alameldeen and D. Wood, "Variability in architectural simulations of multi-threaded workloads," in *International Symposium on High-Performance Computer Architecture (HPCA)*, Feb. 2003, pp. 7–18.
- [3] A. Alameldeen and D. Wood, "IPC considered harmful for multiprocessor workloads," *IEEE Micro*, vol. 26, no. 4, pp. 8–17, 2006.
- [4] D. Albonesi, "Selective cache ways: on-demand cache resource allocation," in *Proceedings of the 32nd Annual ACM/IEEE International Symposium on Microarchitecture*, 1999, pp. 248–259.
- [5] E. K. Ardestani and J. Renau, "ESEC: A fast multicore simulator using time-based sampling," in *International Symposium on High Performance Computer Architecture (HPCA)*, Feb. 2013, pp. 448–459.
- [6] E. Argollo, A. Falcón, P. Faraboschi, M. Monchiero, and D. Ortega, "Cotson: infrastructure for full system simulation," *ACM SIGOPS Operating Systems Review*, vol. 43, no. 1, pp. 52–61, 2009.
- [7] D. Bailey, T. Harris, W. Saphir, R. Van Der Wijngaart, A. Woo, and M. Yarrow, "The NAS parallel benchmarks 2.0," NAS-95-020, NASA Ames Research Center, Tech. Rep., 1995.
- [8] K. C. Barr, H. Pan, M. Zhang, and K. Asanovic, "Accelerating multiprocessor simulation with a memory timestamp record," in *International Symposium on Performance Analysis of Systems and Software (ISPASS)*, Mar. 2005, pp. 66–77.
- [9] C. Bienia, "Benchmarking modern multiprocessors," Ph.D. dissertation, Princeton University, January 2011.
- [10] N. Binkert, B. Beckmann, G. Black, S. K. Reinhardt, A. Saidi, A. Basu, J. Hestness, D. R. Hower, T. Krishna, S. Sardashti, R. Sen, K. Sewell, M. Shoaib, N. Vaish, M. D. Hill, and D. A. Wood, "The gem5 simulator," *SIGARCH Computer Architecture News*, vol. 39, no. 2, pp. 1–7, Aug. 2011.
- [11] W. L. Bircher and L. John, "Predictive power management for multi-core processors," in *Proceedings of the 2010 International Conference on Computer Architecture (ISCA)*, 2010, p. 243–255.
- [12] F. Bodon and L. Rónyai, "Trie: an alternative data structure for data mining algorithms," *Mathematical and Computer Modelling*, vol. 38, no. 7-9, pp. 739–751, 2003.
- [13] J. Bucek, K.-D. Lange, and J. v. Kistowski, "SPEC CPU2017: Next-generation compute benchmark," in *International Conference on Performance Engineering (ICPE)*, Apr. 2018, pp. 41–42.
- [14] A. Butko, R. Garibotti, L. Ost, V. Lapotre, A. Gamatie, G. Sassatelli, and C. Adeniyi-Jones, "A trace-driven approach for fast and accurate simulation of manycore architectures," in *The 20th Asia and South Pacific Design Automation Conference (ASP-DAC)*, 2015, pp. 707–712.
- [15] T. E. Carlson, W. Heirman, and L. Eeckhout, "Sniper: Exploring the level of abstraction for scalable and accurate parallel multi-core simulation," in *International Conference for High Performance Computing, Networking, Storage and Analysis (SC)*, Nov. 2011, pp. 52:1–52:12.
- [16] T. E. Carlson, W. Heirman, and L. Eeckhout, "Sampled simulation of multi-threaded applications," in *International Symposium on Performance Analysis of Systems and Software (ISPASS)*, Apr. 2013, pp. 2–12.
- [17] T. E. Carlson, W. Heirman, K. Van Craeynest, and L. Eeckhout, "BarrierPoint: Sampled simulation of multi-threaded applications," in *International Symposium on Performance Analysis of Systems and Software (ISPASS)*, Mar. 2014, pp. 2–12.
- [18] M. Casas, H. Servat, R. M. Badia, and J. Labarta, "Extracting the optimal sampling frequency of applications spectral analysis," *Concurrency and Computation: Practice and Experience*, vol. 24, pp. 237–259, 03 2012.
- [19] J. Charles, P. Jassi, N. S. Ananth, A. Sadat, and A. Fedorova, "Evaluation of the intel® core™ i7 turbo boost feature," in *2009 IEEE International Symposium on Workload Characterization (IISWC)*, 2009, pp. 188–197.
- [20] D. Chasapis, M. Casas, M. Moretó, R. Vidal, E. Ayguadé, J. Labarta, and M. Valero, "Parsecs: Evaluating the impact of task parallelism in the parsec benchmark suite," *ACM Transactions on Architecture and Code Optimization (TACO)*, vol. 12, no. 4, pp. 1–22, 2015.
- [21] T. Chen, Y. Chen, Q. Guo, O. Temam, Y. Wu, and W. Hu, "Statistical performance comparisons of computers," in *IEEE International Symposium on High-Performance Computer Architecture (HPCA)*, 2012, pp. 1–12.
- [22] W.-K. Chen, S. Lerner, R. Chaiken, and D. M. Gillies, "Mojo: A dynamic optimization system," in *3rd ACM Workshop on Feedback-Directed and Dynamic Optimization (FDDO-3)*, 2000, pp. 81–90.
- [23] S. Dasgupta, "Experiments with random projection," in *Proceedings of the Sixteenth Conference on Uncertainty in Artificial Intelligence*, 2000, pp. 143–151.
- [24] S. De Pestel, S. Van den Steen, S. Akram, and L. Eeckhout, "RPPM: Rapid performance prediction of multithreaded workloads on multicore processors," in *International Symposium on Performance Analysis of Systems and Software (ISPASS)*, Mar. 2019, pp. 257–267.
- [25] S. De Pestel, S. Eyerman, and L. Eeckhout, "Micro-architecture independent branch behavior characterization," in *2015 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS)*, 2015, pp. 135–144.
- [26] M. Deilmann et al., "A guide to vectorization with Intel C++ compilers," *Intel Corporation*, pp. 20–21, 2012.
- [27] A. Diavastos and P. Trancoso, "Switches: A lightweight runtime for dataflow execution of tasks on many-cores," *ACM Transactions on Architecture and Code Optimization (TACO)*, vol. 14, no. 3, pp. 1–23, 2017.
- [28] J. Doweck, W.-F. Kao, A. K.-y. Lu, J. Mandelblat, A. Rahatekar, L. Rappoport, E. Rotem, A. Yasin, and A. Yoaz, "Inside 6th-generation Intel Core: New microarchitecture code-named Skylake," *IEEE Micro*, vol. 37, no. 2, pp. 52–62, 2017.
- [29] A. Duran, E. Ayguadé, R. M. Badia, J. Labarta, L. Martinell, X. Martorell, and J. Planas, "OmpSs: a proposal for programming heterogeneous multi-core architectures," *Parallel processing letters*, vol. 21, no. 02, pp. 173–193, 2011.
- [30] L. Eeckhout, Y. Luo, K. De Bosschere, and L. K. John, "BLRL: Accurate and efficient warmup for sampled processor simulation," *The Computer Journal*, vol. 48, no. 4, pp. 451–459, 2005.
- [31] D. Eklov and E. Hagersten, "StatStack: Efficient modeling of LRU caches," in *2010 IEEE International Symposium on Performance Analysis of Systems & Software (ISPASS)*, 2010, pp. 55–65.
- [32] C. Elkan, "Using the triangle inequality to accelerate k-means," in *Proceedings of the 20th International Conference on Machine Learning (ICML)*, 2003, pp. 147–153.
- [33] S. Eyerman and L. Eeckhout, "Fine-grained DVFS using on-chip regulators," *ACM Transactions on Architecture and Code Optimization (TACO)*, vol. 8, no. 1, pp. 1–24, 2011.
- [34] S. Eyerman, L. Eeckhout, T. Karkhanis, and J. E. Smith, "A mechanistic performance model for superscalar out-of-order processors," *ACM*

Transactions on Computer Systems (TOCS), vol. 27, no. 2, pp. 1–37, 2009.

- [35] S. Eyerman, W. Heirman, S. Van den Steen, and I. Hur, “Enabling branch-mispredict level parallelism by selectively flushing instructions,” in *54th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*, 2021, pp. 767–778.
- [36] A. Falcón, P. Faraboschi, and D. Ortega, “Combining simulation and virtualization through dynamic sampling,” in *2007 IEEE International Symposium on Performance Analysis of Systems & Software (ISPASS)*, 2007, pp. 72–83.
- [37] E. W. Forgy, “Cluster analysis of multivariate data: efficiency versus interpretability of classifications,” *Biometrics*, vol. 21, pp. 768–769, 1965.
- [38] T. Grass, A. Rico, M. Casas, M. Moreto, and E. Ayguadé, “TaskPoint: Sampled simulation of task-based programs,” in *International Symposium on Performance Analysis of Systems and Software (ISPASS)*, Apr. 2016, pp. 296–306.
- [39] A. Hajiabadi, A. Diavastos, and T. E. Carlson, “NOREBA: A compiler-informed non-speculative out-of-order commit processor,” in *Proceedings of the 26th ACM International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, 2021, p. 182–193.
- [40] J. W. Haskins and K. Skadron, “Memory reference reuse latency: Accelerated warmup for sampled microarchitecture simulation,” in *IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS)*, 2003, pp. 195–203.
- [41] S. Hassani, G. Southern, and J. Renau, “LiveSim: Going live with microarchitecture simulation,” in *International Symposium on High Performance Computer Architecture (HPCA)*, Mar. 2016, pp. 606–617.
- [42] C. Isci, A. Buyuktosunoglu, C.-Y. Cher, P. Bose, and M. Martonosi, “An analysis of efficient multi-core global power management policies: Maximizing performance for a given power budget,” in *2006 39th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*, 2006, pp. 347–358.
- [43] C. Jiang, Z. Yu, H. Jin, C. Xu, L. Eeckhout, W. Heirman, T. E. Carlson, and X. Liao, “PCantorSim: Accelerating parallel architecture simulation through fractal-based sampling,” *ACM Transactions on Architecture and Code Optimization (TACO)*, vol. 10, no. 4, Dec. 2013.
- [44] J. L. Kihm, S. D. Strom, and D. A. Connors, “Phase-guided small-sample simulation,” in *2007 IEEE International Symposium on Performance Analysis of Systems & Software (ISPASS)*, 2007, pp. 84–93.
- [45] W. Kim, M. S. Gupta, G.-Y. Wei, and D. Brooks, “System level analysis of fast, per-core dvfs using on-chip switching regulators,” in *2008 IEEE 14th International Symposium on High Performance Computer Architecture (HPCA)*, 2008, pp. 123–134.
- [46] N. Kulkarni, F. Qi, and C. Delimitrou, “Pliant: Leveraging approximation to improve datacenter resource efficiency,” in *2019 IEEE International Symposium on High Performance Computer Architecture (HPCA)*, 2019, pp. 159–171.
- [47] J. Lau, J. Sampson, E. Perelman, G. Hamerly, and B. Calder, “The strong correlation between code signatures and performance,” in *IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS)*, 2005.
- [48] J. Lau, E. Perelman, and B. Calder, “Selecting software phase markers with code structure analysis,” in *International Symposium on Code Generation and Optimization (CGO)*, Mar. 2006, pp. 135–146.
- [49] “LoopPoint tools,” <https://github.com/nus-comparch/looppoint>.
- [50] R. Mattson, J. Gecsei, D. R. Slutz, and I. L. Traiger, “Evaluation techniques for storage hierarchies,” *IBM Systems Journal*, vol. 9, no. 2, pp. 78–117, 1970.
- [51] S. Mittal, Y. Cao, and Z. Zhang, “Master: A multicore cache energy-saving technique using dynamic cache reconfiguration,” *IEEE Transactions on very large scale integration (VLSI) systems*, vol. 22, no. 8, pp. 1653–1665, 2013.
- [52] S. Mittal, Z. Zhang, and J. S. Vetter, “Flexiway: A cache energy saving technique using fine-grained cache reconfiguration,” in *2013 IEEE 31st International Conference on Computer Design (ICCD)*, 2013, pp. 100–107.
- [53] N. Nikoleris, A. Sandberg, E. Hagersten, and T. E. Carlson, “CoolSim: Statistical techniques to replace cache warming with efficient, virtualized profiling,” in *2016 International Conference on Embedded Computer*
- [54] N. Nikoleris, L. Eeckhout, E. Hagersten, and T. E. Carlson, “Directed statistical warming through time traveling,” in *International Symposium on Microarchitecture (MICRO)*, Oct. 2019, pp. 1037–1049.
- [55] “OpenMP 3.1 API C/C++ Syntax Quick Reference Card,” <https://www.openmp.org/wp-content/uploads/OpenMP3.1-CCard.pdf>.
- [56] H. Patil and T. E. Carlson, “Pinballs: portable and shareable user-level checkpoints for reproducible analysis and simulation,” in *Workshop on Reproducible Research Methodologies (REPRODUCE)*, Feb. 2014.
- [57] H. Patil, A. Isaev, W. Heirman, A. Sabu, A. Hajiabadi, and T. E. Carlson, “ELFies: Executable region checkpoints for performance analysis and simulation,” in *International Symposium on Code Generation and Optimization (CGO)*, Feb./Mar. 2021, pp. 126–136.
- [58] H. Patil, C. Pereira, M. Stallcup, G. Lueck, and J. Cownie, “PinPlay: A framework for deterministic replay and reproducible analysis of parallel programs,” in *International Symposium on Code Generation and Optimization (CGO)*, Apr. 2010, pp. 2–11.
- [59] E. Perelman, M. Polito, J.-Y. Bouguet, J. Sampson, B. Calder, and C. Dulong, “Detecting phases in parallel applications on shared memory architectures,” in *International Parallel Distributed Processing Symposium (IPDPS)*, Apr. 2006.
- [60] A. Sabu, H. Patil, W. Heirman, and T. E. Carlson, “LoopPoint: Checkpoint-driven sampled simulation for multi-threaded applications,” in *International Symposium on High Performance Computer Architecture (HPCA)*, 2022.
- [61] D. Sanchez and C. Kozyrakis, “ZSim: Fast and accurate microarchitectural simulation of thousand-core systems,” in *International Symposium on Computer Architecture (ISCA)*, Jun. 2013, pp. 475–486.
- [62] A. Sandberg, N. Nikoleris, T. E. Carlson, E. Hagersten, S. Kaxiras, and D. Black-Schaffer, “Full Speed Ahead: Detailed architectural simulation at near-native speed,” in *2015 IEEE International Symposium on Workload Characterization (IISWC)*, 2015, pp. 183–192.
- [63] X. Shen, Y. Zhong, and C. Ding, “Locality phase prediction,” in *Proceedings of the 11th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, 2004, p. 165–176.
- [64] T. Sherwood, E. Perelman, G. Hamerly, and B. Calder, “Automatically characterizing large scale program behavior,” in *International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, Oct. 2002, pp. 45–57.
- [65] T. Sherwood, S. Sair, and B. Calder, “Phase tracking and prediction,” in *Proceedings of the 30th Annual International Symposium on Computer Architecture (ISCA)*, 2003, p. 336–349.
- [66] M. Van Biesbrouck, B. Calder, and L. Eeckhout, “Efficient sampling startup for simpoint,” *IEEE Micro*, vol. 26, no. 4, pp. 32–42, 2006.
- [67] M. J. Voss and R. Eigemann, “High-level adaptive program optimization with adapt,” in *Proceedings of the Eighth ACM SIGPLAN Symposium on Principles and Practices of Parallel Programming (PPoPP)*, 2001, pp. 93–102.
- [68] T. F. Wenisch, R. E. Wunderlich, M. Ferdman, A. Ailamaki, B. Falsafi, and J. C. Hoe, “SimFlex: Statistical sampling of computer system simulation,” *IEEE Micro*, vol. 26, no. 4, pp. 18–31, 2006.
- [69] R. E. Wunderlich, T. F. Wenisch, B. Falsafi, and J. C. Hoe, “SMARTS: Accelerating microarchitecture simulation via rigorous statistical sampling,” in *International Symposium on Computer Architecture (ISCA)*, Jun. 2003, pp. 84–97.
- [70] X. You, C. Liu, H. Yang, P. Wang, Z. Luan, and D. Qian, “Vectorizing SpMV by exploiting dynamic regular patterns,” in *Proceedings of the 51st International Conference on Parallel Processing (ICPP)*, 2022, pp. 1–12.
- [71] C. Yount, H. Patil, and M. S. Islam, “Graph-matching-based simulation-region selection for multiple binaries,” in *International Symposium on Performance Analysis of Systems and Software (ISPASS)*, Mar. 2015, pp. 52–61.
- [72] J. Zeng, H. Kim, J. Lee, and C. Jung, “Turnpike: Lightweight soft error resilience for in-order cores,” in *54th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*, 2021, p. 654–666.