

SmartOClock: Workload- and Risk-Aware Overclocking in the Cloud

Jovan Stojkovic², Pulkit A. Misra¹, Íñigo Goiri¹, Sam Whitlock¹, Esha Choukse¹,
Mayukh Das¹, Chetan Bansal¹, Jason Lee¹, Zoey Sun¹, Haoran Qiu²
Reed Zimmermann³, Savyasachi Samal¹, Brijesh Warriar¹, Ashish Raniwala¹, Ricardo Bianchini¹

¹Microsoft

²University of Illinois at Urbana Champaign

³University of Texas at Austin

Abstract—Operating server components beyond their voltage and power design limit (*i.e.*, overclocking) enables improving performance and lowering cost for cloud workloads. However, overclocking can significantly degrade component lifetime, increase power draw, and cause power capping events, eventually diminishing the performance benefits.

In this paper, we characterize the impact of overclocking on cloud workloads by studying their profiles from production deployments. Based on the characterization insights, we propose SmartOClock, the first distributed overclocking management platform specifically designed for cloud environments. SmartOClock is a workload-aware scheme that relies on power predictions to heterogeneously distribute the power budgets across its servers based on their needs and then enforce budget compliance locally, per-server, in a decentralized manner.

SmartOClock reduces the tail latency by 9%, application cost by 30% and total energy consumption by 10% for latency-sensitive microservices on a 36-server deployment. Simulation analysis using production traces show that SmartOClock reduces the number of power capping events by up to 95% while increasing the overclocking success rate by up to 62%. We also describe lessons from building a first-of-its-kind overclockable cluster in Microsoft Azure for production experiments.

I. INTRODUCTION

Motivation. Cloud services provision resources to meet their peak performance requirements [21], [25], [39], [62], [81]. For example, many services need to keep their high-percentile latency (*e.g.*, P99) below a predetermined Service-Level Objective (SLO) [24]. These services incur high operating costs to reserve enough resources for handling infrequent load spikes and leave a substantial portion underutilized or even idle for the majority of time when their load is below its peak.

As an example, **Figure 1** illustrates the aggregate load pattern on a typical weekday of three services that are part of Microsoft’s productivity and collaboration suite. Collectively, these three services use $\sim 1\text{M}$ virtual cores (across regions) to handle peaks that last for a few hours per day - between 10 am to noon for *Service A* and 5 minutes at the top and bottom of the hour for the other two services.

Emerging cloud paradigms, such as autoscaling [4], [33], [71] and serverless computing [5], [35], [45], [70], can be used to dynamically remove and add Virtual Machine (VM) instances for managing cost. However, these solutions (1) can increase the application’s tail latency as booting up a new

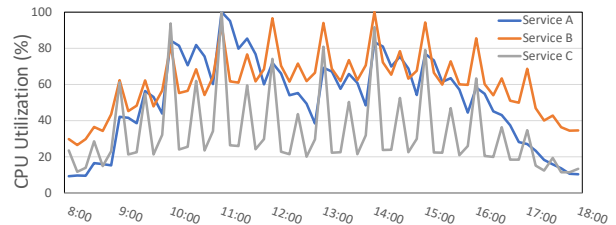


Fig. 1: Load pattern on a typical weekday in one region. Utilization normalized to the peak of each service.

VM can take up to a few minutes [1], and (2) cannot be easily applied for stateful services [46], [69]. Hence, many applications still statically provision for infrequent load spikes.

On the other hand, advances in processing and datacenter cooling technologies have enabled component (*e.g.*, CPU, GPU) overclocking, *i.e.*, operation beyond typical voltage and power design limit [51]. Overclocking boosts a workload’s performance and enables handling transient load spikes in a cost-efficient manner. For example, CPU overclocking during a service’s peak can keep the tail latency below the required SLO, while saving cost by reducing provisioned resources.

However, overclocking is not free. If used naively, it increases power draw and can cause frequent power capping events that diminish performance benefits. Worse, it degrades component lifetime (or reliability) through accelerated wear-out and, thus, cannot be used indefinitely. The limited amount of overclocking needs to be used smartly as it may not benefit all workloads at all times: (1) overclocking the CPU of a memory-bound workload, or (2) overclocking a workload while experiencing a low load will not provide much benefit. Finally, providers also need to protect workload SLOs when overclocking is unavailable. For example, a workload might have under-provisioned due to reliance on overclocking, but it would miss its SLOs under peak load if its VMs cannot be overclocked. Therefore, providers must use overclocking carefully while managing the associated risks.

Our work. For efficient use of overclocking in the cloud, we analyze cloud workloads and production traces, including the services from **Figure 1**. We observe the following. First, overclocking improves the performance of popular cloud workloads. However, a workload-agnostic overclocking scheme

is suboptimal and often leads to missed SLOs or wasted overclocking cycles. Second, power and lifetime headroom exists to overclock most of the times without triggering power capping or compromising on reliability. Third, resource utilization history can be used to predict the availability of power and reliability impact from overclocking. Fourth, servers’ power draw within a power delivery unit (*e.g.*, a rack) is diverse, but the limit is still evenly distributed which disproportionately hurts performance of power-hungry servers during a capping event. However, predictability in power draw enables assigning heterogeneous limits. Finally, a decentralized approach for power draw enforcement enables servers to find an efficient limit in case of initial assignment mispredictions.

We use the characterization insights to design SmartOClock, the first distributed overclocking management platform for the cloud. It enables a wide variety of cloud workloads to run with high performance at a lower cost. SmartOClock achieves its goals through four novel design principles.

First, SmartOClock uses bidirectional communication with the application to maximize the application’s benefits from overclocking. Applications can use metrics (*e.g.*, latency, CPU utilization) or schedule-based policies and the overclocking decisions can be made based on instance- and deployment-level monitoring. Second, it uses *admission control* to reserve power (from any headroom) and overclocking budget for workloads. This step provides a predictable overclocking experience for workloads and SmartOClock can take corrective actions, like scale-out, if it is unable to honor a reservation. Third, it leverages power predictability for assigning *heterogeneous* server power budgets, which provide better performance during capping for power safety. Finally, SmartOClock makes *decentralized* overclocking decisions for improved fault tolerance. Each server takes local decisions for granting overclocking requests based on its assigned power and overclocking budgets. It can also perform explorations to revise inefficient assignments (*e.g.*, due to mispredictions).

We evaluate SmartOClock on a real server cluster and through simulations by using production traces. The cluster evaluation is performed on 36 overclockable servers (across 2-racks) running latency-sensitive microservices as candidates for overclocking and throughput-optimized power hungry machine learning (ML) training workloads, which are not overclocked. Our results show that SmartOClock reduces the P99 latency by 8.9% and application cost by 30.4% for latency-sensitive microservices, and total cluster energy consumption by 10% over state-of-the-art autoscaling solution. To validate our findings at scale, we use traces from hundreds of production racks and simulate SmartOClock. When compared to all practical policies, SmartOClock reduces the number of power capping events by up to 94.7% while increasing the overclocking success rate by up to 61.8%. We have also created a 2-rack overclockable cluster for production experimentation and share some lessons in [Section VI](#).

Related work. While there is a rich body of work on CPU turbo-boost [16], [18], [30], [55], [74], [79], [101] and

datacenter power management [37], [57], [59], [80], [84], [94], [97], overclocking introduces unique challenges not addressed by the prior work. First, a cloud provider does not need to manage any reliability impact from turbo since CPU vendors design turbo to meet a provider’s lifetime requirements. Cloud CPUs operate in performance mode, which always operates them at the highest turbo frequency within constraints (*e.g.*, power, thermal) [18], [82]. Vendors do not specify turbo timing limitations nor advise software-level core wear leveling in their warranty terms [8], [47], and non-judicious turbo use does not degrade reliability [76]. Generally, CPU failure is amongst the lowest types of failure in cloud servers [65], [91]. Second, the power oversubscription policies factor the higher demand from turbo. Although this approach increases the total cost of ownership, it is necessary to meet the performance Service-Level Agreements (SLAs) [7], [36], [72]. In contrast, overclocking (beyond turbo) further improves performance but has a reliability impact that is not covered at design time by the vendors. Furthermore, a provider does not need to provision power for overclocking since turbo is sufficient to meet its performance SLAs. Therefore, overclocking is opportunistic - a provider needs to manage the power and reliability impact, while protecting workload SLOs when overclocking is unavailable; a problem setting not explored by prior work.

Summary. We make the following main contributions:

- We characterize the opportunities and challenges of overclocking cloud workloads, including the impact on power and component lifetime.
- We propose SmartOClock, a distributed overclocking management platform specifically designed for the cloud.
- We evaluate SmartOClock in a real system running latency-critical workloads, and using large-scale production traces.
- We share lessons from overclocking production workloads.

II. BACKGROUND

Power management in cloud datacenters. The power delivery system in a cloud datacenter is organized in a hierarchy [57], [84], [94], [97]; the power budget of each parent node is split equally among its children. As providers oversubscribe power to improve utilization, the sum of the peak power draw of children nodes can exceed the budget of the parent (*e.g.*, servers in a rack) [57], [84], [94]. Under normal operation, child nodes can draw more than their even share if the cumulative power is below the parent’s limit. When it exceeds a threshold, power capping mechanisms (*e.g.*, Intel RAPL [22], prioritized capping [57], [59]) are used for safety. These mechanisms hurt performance as they reduce CPU frequency and can even throttle memory to restrict server power. To meet their performance SLAs, providers carefully oversubscribe to minimize/avoid capping events.

Component overclocking. Prior work shows the feasibility of overclocking in the cloud [51]. Overclocking operates components (*e.g.*, CPUs, GPUs) beyond their specifications to get frequencies even higher than turbo [10], [50].

A large fraction of cloud workloads, such as search or video conferencing [21], [89], are user-facing applications with transient load spikes. These workloads collectively consume millions of virtual cores to handle peak load. For Microsoft’s productivity and collaboration services, although chat and conference calls occur throughout the day, the peak that governs resource provisioning lasts for a few hours each day (Figure 1). Overclocking can be used during these peaks to save costs. However, a provider needs to manage the risks from overclocking. For example, for reliability management, the peak duration needs to be within the daily overclocking budget (*e.g.*, 10% per day) that satisfies component lifetime goals. Overclocking impacts reliability [3] due to three main reasons: (1) gate oxide breakdown, (2) electro-migration, and (3) thermal cycling. These processes are time-dependent and accelerate the lifetime reduction. Prior work has showed that there is an exponential relationship between temperature, voltage, and component lifetime [27], [51], [66], [93], [96].

III. CHALLENGES AND OPPORTUNITIES

A successful overclocking management scheme needs to satisfy workload performance requirements, while managing the impact of overclocking on power and component lifetime. To design such a scheme, we answer the following questions.

Q1: When do workloads benefit from overclocking? To efficiently use overclocking, a cloud platform needs to understand workloads’ behavior and needs. Treating VMs as opaque and using workload performance proxies (*e.g.*, instructions per cycle (IPC), CPU utilization) for overclocking can be suboptimal as the relationship between proxies and target performance metric is not always clear. Without knowing a workload’s performance goals, the platform may overclock prematurely (*i.e.*, under low load that does not impact tail performance) and, due to the lifetime impact, lose the ability to overclock when really needed. Combining IPC with CPU utilization as a proxy for load can be inefficient too because the performance of some workloads is impacted at a moderate CPU utilization while others are unimpacted even under high utilization. Finally, operators can even have *deployment-level* goals for provisioning (number of VMs) and overclocking based on instance-level monitoring only will be inefficient.

To illustrate these scenarios, we profile two classes of popular cloud workloads: (1) microservices from the largest open-source benchmark suite, DeathStarBench [32], and (2) a proprietary web conferencing application called WebConf.

Microservices. We run eight SocialNet microservices [32] under varying loads (low, medium, and high) in three environments: *Baseline*, *Overclock*, and *ScaleOut*. *Baseline* and *Overclock* run a single VM at turbo (3.3 GHz) and overclocked (4.0 GHz) frequency. *ScaleOut* has two VMs running at turbo. Figure 2 shows the tail latency of the microservices. The red horizontal line indicates SLO, where the SLO for each service is set to be 5 times its execution time on an unloaded system [26], [60], [73]. Figure 3 shows their CPU utilization.

ScaleOut is provisioned to handle the peak load and always operates 2 VMs that run at turbo. Although it provides the

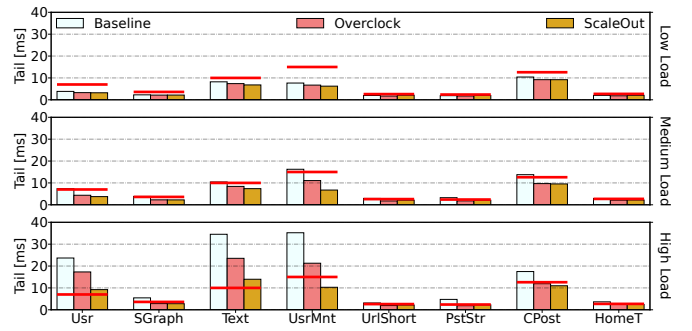


Fig. 2: Tail latency of SocialNet microservices with different loads in Baseline, Overclock, and ScaleOut environments.

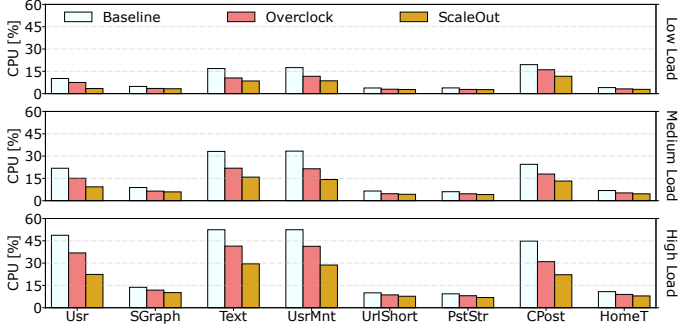


Fig. 3: CPU utilization of SocialNet microservices with different loads in Baseline, Overclock, and ScaleOut environments.

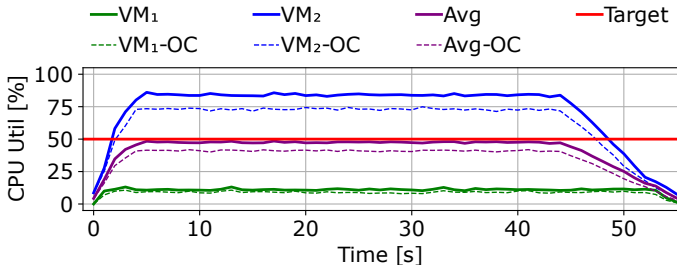


Fig. 4: CPU utilization timeline with and without overclocking for two WebConf VMs.

best performance, it also incurs the highest cost. In contrast, Overclock uses a single VM and still keeps the tail latency below the SLO in many cases, thereby avoiding the need to scale out. However, some services (*e.g.*, Usr) can tolerate higher CPU utilization without violating their SLO while others (*e.g.*, UriShort) violate their SLO even under a low utilization. Therefore, a workload-agnostic policy using CPU utilization for overclocking will make suboptimal decisions. These observations hold for any cloud workload with similar characteristics – bursty load with tail latency as the key metric. For example, ML inference servers [60], [98], serverless computing [86], and key-value stores [61] amongst others.

WebConf. The workload hosts conferences in a VM. For fault-tolerance, operators provision VMs across availability zones (AZ) in a region. In an AZ, provisioning keeps the average *deployment-level* CPU utilization below 50% to handle load from another failed AZ. Overclocking can save cost for

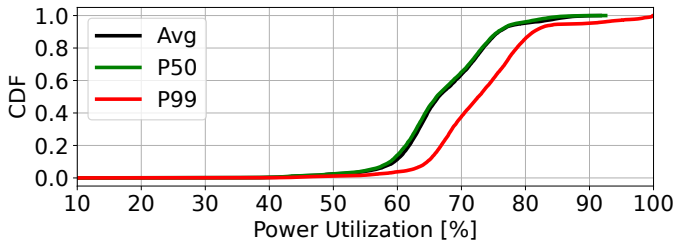


Fig. 5: Average, median (P50), and peak (P99) power utilization of 7,100 racks over 6 weeks in three regions.

WebConf through deployment-level decisions. Individual VMs can have high utilization, but overlocking them is suboptimal since the deployment-level utilization may be below the target.

To illustrate, we execute WebConf on two VMs. VM_1 has a low load while VM_2 's load is high. Figure 4 shows the VM- and deployment-level average CPU utilization. Although overlocking provides a benefit, it is unnecessary since the baseline already meets the workload performance (provisioning) goal.

Q2: Are there enough resources for overlocking? Since overlocking increases power draw and component wear out, we need headroom for these resources.

Power headroom. We analyze the power draw of 7.1k dedicated racks that run Microsoft's productivity and collaboration services, including those from Figure 1, which are used by millions of users across the world. The racks span all major regions (e.g., United States, Europe, Asia) and each rack has 24-32 servers. The analysis period is 6 weeks (April 10th – May 12th, 2023). Figure 5 shows the CDF of average, median (P50), and P99 rack power utilization. Half the racks have an average utilization lower than 66%. Importantly, 50% and 90% of the racks have P99 lower than 73% and 89%, respectively. We observe similar power patterns on non-dedicated racks in Azure with a mix of first- and third-party workloads.

To estimate the power impact from overlocking, we use the overlocking requirements of critical user-facing workloads that constitute 45% of the deployed cores. Their requirements vary – some require overlocking for several minutes per hour, while others for multiple hours per weekday. Figure 6 shows the power draw of a rack without and with overlocking for five weekdays; the red line shows the rack power limit. Each server in this rack hosts VMs of many distinct services and captures a typical datacenter environment with multi-tenant servers. The rack power draw is below the limit for the baseline, but overlocking exceeds the limit and causes capping. More generally, overlocking the selected workloads will not result in capping for 85% of the time. For the remaining 15%, naive overlocking causes 30-50% degradation in workload performance (core frequency) due to capping. However, there is still headroom available on these power-constrained racks, but it is insufficient to overlock to the highest frequency; the available headroom is 75% of the requisite at P99.

Therefore, most of the time (85%) racks have the needed power headroom for overlocking. However, a power-aware policy is needed for the constrained scenarios.

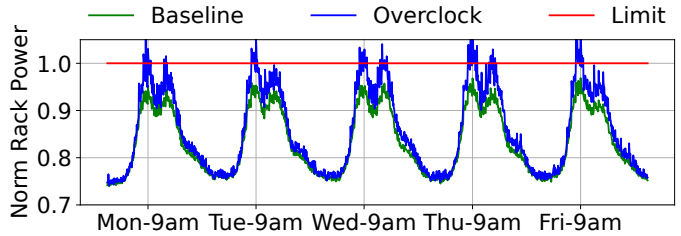


Fig. 6: Example of rack power draw over 5 weekdays

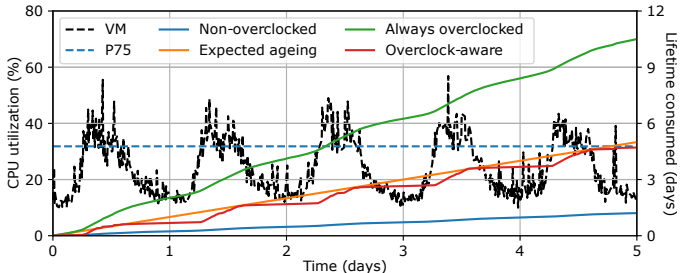


Fig. 7: CPU ageing for a VM running a workload with a diurnal pattern under multiple overlocking policies.

These findings are with the default Azure VM scheduler that uses a set of resource-centric placement rules [40]. Providers can add power-aware scheduling policies to aid overlocking, but this exploration is future work. Nonetheless, even with optimized placement, there will still be power-constrained scenarios where overlocking has to be performed carefully.

Component lifetime headroom. Prior work shows that advanced cooling (e.g., wax, immersion) is needed for enabling sprinting/overlocking [30], [51], [78], [79]. However, there is opportunity to overlock even in air-cooled server deployments. Cloud server cooling is designed for operating components at their rated thermal design power (TDP). However, servers rarely consume their TDP due to low resource utilization in the cloud [21]. Several factors contribute to the low utilization. First, over-provisioning and diurnal workload patterns result in low VM utilization. Second, workload heterogeneity on servers results in low server utilization. Each server hosts many small VMs (2-8 cores). For resiliency, operators spread their VMs across servers and racks. Consequently, the VMs on any given server belong to different workloads. This heterogeneity results in low server utilization as the workloads have different peak times. Consequently, components are not thermally constrained for overlocking in air and advanced cooling can be used to enhance the capability (e.g., duration) as lower operating temperatures reduce ageing [51]. Finally, since overlocking does not exceed the TDP nor the rack limit, it will not cause additional cooling-related failures.

In fact, under-utilization enables overlocking in air. Vendors assume near-100% usage for determining frequencies/voltage (e.g., turbo) that satisfy the lifetime goals. Under-utilization accumulates lifetime credits that can be consumed via overlocking. To understand the opportunity, we use a 7nm composite processor model from TSMC. It uses a complex relation between overlocking (voltage scaling) and CPU utilization (time period at the specified voltage) to model the

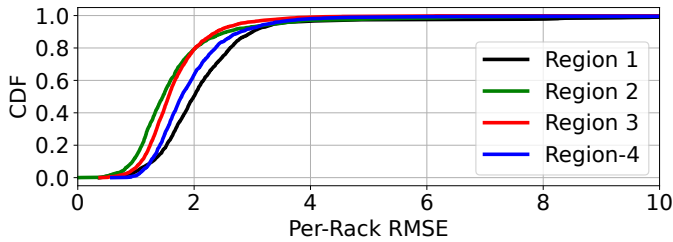


Fig. 8: CDF of RMSE with the power draw patterns of our predictions across 7.1k racks in four regions.

ageing from wear-out in the form of gate oxide failure [23], [58]. The model predicts that a CPU ages by 2.5 years over a 5-year period for a conservative fleet usage. The remaining 2.5 years can be used for overclocking. But naively overclocking for 50% of the time ages the CPU by 5 years in less than a year use due to accelerated wearout. A smart system can constrain overclocking so that the part ages according to the reference (*i.e.*, 1 year ageing over a 1-year period).

Figure 7 illustrates the effect of overclocking policies on ageing. It shows the 5-day CPU utilization of a production workload with a diurnal pattern of daily midday peaks above 50% and valleys lower than 20% at night. The expectation is that the processor ages 5 days over the same period (“Expected ageing”). However, the actual ageing is less than 2-days for the “Non-overclocked” baseline. “Always overclock” ages the CPU over 10 days, indicating that, for the same CPU utilization, overclocking significantly increases wearout. On the other hand, an “Overclock-aware” policy can consume the accumulated credits by overclocking for 25% of the time and not exceed the expected ageing. Offline modeling assumes CPU utilization is unchanged while overclocking for worst-case analysis. However, overclocking’s ageing impact will be less if the utilization reduces. To address this limitation, we are working with the CPU vendors on “wearout counters” for online calculation of the ageing impact (see §VI).

Therefore, overclocking is enabled due to under-utilization and can be improved with advanced cooling. A system must carefully manage overclocking to comply with lifetime goals.

Q3: Can we predict the availability of resources? An efficient overclocking system must perform admission control based on available power and lifetime. We observe that a prediction-based approach can yield high accuracy.

Power predictability. A system needs to predict how much power can be used by overclocking without triggering capping. Figure 6 shows the baseline power draw of a rack and gives us insights that historical observations of power profiles can be leveraged for prediction. The rack hosts multiple services, where each service can have a distinct power profile. However, due to statistical multiplexing, the combined power draw of the rack with heterogeneous services shows a repeatable pattern. We analyzed the power predictability of 7.1K racks (thousands of servers) that collectively run >100 services. Although the racks are dedicated for Microsoft’s productivity and collaboration services, this dataset accounts the fact that

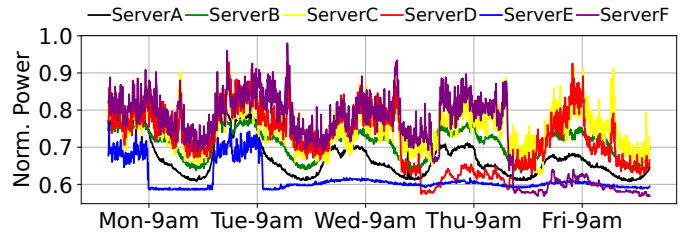


Fig. 9: Normalized power draw over time of six randomly chosen servers within the same rack.

racks and servers on a public cloud host heterogeneous workloads. Furthermore, the dynamicity of cloud platforms (*e.g.*, VM churn according to a workload’s needs) is also reflected.

Figure 8 shows the CDF of Root Mean Squared Error (RMSE) of rack power predictions of four Azure regions. The RMSE is low even at high percentiles indicating high predictability. For example, in Region 3, 50% and 99% of the racks have an RMSE lower than 1.95W and 5.11W, respectively. The findings are similar in the other regions. Furthermore, an analysis of 20K non-dedicated racks, running a mix of first- and third-party workloads, in the three most popular Azure regions yielded similar results. A major reason for this predictability is long-lived VMs that govern resource utilization. Prior work shows that long-lived VMs (or jobs) account for >95% of allocated resources [21], [81], [85].

Component lifetime impact predictability. To remain within the overclocking lifetime budget, a system needs to predict how much overclocked CPU cycles a given workload will consume. As a server’s power draw depends on CPU utilization, predictability in power indicates predictability in CPU utilization. Using the aforementioned methodology for a rack’s power, we now analyze the CPU utilization predictability. Our results show that CPU utilization of servers are also predictable: more than 50% and 90% of the servers have an RMSE of CPU utilization lower than 3.13% and 7.82%, respectively.

Therefore, historical observations of power draw and CPU utilization can be used to predict the available power and component lifetime headroom for overclocking.

Q4: How to assign power budgets? A server’s power budget for “safe” overclocking depends on the power draw of the other servers in the hierarchy (*e.g.*, a rack). Under fair share, the rack power budget is split equally across all servers and each server can locally ensure that its power draw stays below the limit to avoid capping while overclocking. However, this approach is inefficient since some servers may not be able to overclock even while the rack is not power-constrained.

Figure 9 shows the normalized power draw over 5 weekdays of six randomly chosen servers in a rack; each server is a different color. We can see that servers have very different power profiles. Some servers may use even 30% less power than others. In addition, servers that consume the most power in a rack change over time. For example, at different timestamps, ServerC, ServerD, or ServerF may be the power dominant one.

Therefore, an efficient overclocking system needs to split the rack power budget *heterogeneously* across servers. His-

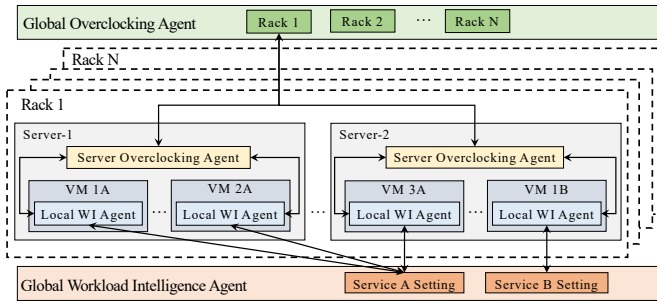


Fig. 10: Overview of the SmartOClock overlocking system.

torical observations of server power demand and rack-level headroom can be used for the heterogeneous attribution.

Q5: How to efficiently use the power? The power headroom for overlocking in a rack is consumed by all servers in that rack. Thus, to grant or reject an overlocking request, each server should contact a centralized entity that has the global view of the rack’s total remaining power headroom. Unfortunately, this approach is expensive and limits the system’s fault-tolerance – if the centralized entity fails, then all overlocking requests would be rejected. Making local overlocking decisions using assigned server power budgets improves fault tolerance. However, overlocking requests may still be rejected due to inefficient assignments. For example, a scheme that uses power predictions for budget assignments can be suboptimal due to mispredictions.

Therefore, a high-performance and fault-tolerant overlocking system needs to be decentralized and should allow servers to explore beyond their potentially stale power limits.

IV. SMARTOLOCK

Driven by the characterization insights, we propose *SmartOClock*: a distributed overlocking management platform for the cloud. It is readily integrated with existing platforms and enables a wide variety of workloads to run with high performance at lower cost. SmartOClock responds to the outlined questions for an efficient overlocking scheme through four novel features. First, it is *workload-intelligent* as it uses hints provided by workloads to extract the most benefits from overlocking. Second, SmartOClock performs *prediction-based admission control* of overlocking requests to avoid power capping and premature component wearout. Third, it uses predictions to split the rack power limit *heterogeneously* across servers. Finally, SmartOClock uses a *decentralized scheme* for budget enforcement while overlocking and allows controlled exploration to revise inefficient assignments.

Architecture. Figure 10 shows the architecture of SmartOClock. The system is organized hierarchically where each controller manages the components on its level and communicates with the controllers from the upper and lower levels. First, when deploying their services, the workload owners configure the *Global Workload Intelligence Agent* for their service. They specify the conditions under which the workload needs to be overlocked. As workloads are composed of one or more

VMs, each VM is deployed with its own *Local Workload Intelligence Agent*. Like conventional auto-scaling, the local agent collects the metrics of interest from the VM and sends them to the global agent. Thus, this setup does not introduce new security or privacy challenges. The global agent uses the metrics to decide if any VM needs to be overlocked and sends a signal to the local agent of such VMs. On receiving a signal, a local agent sends an overlocking request to the *Server Overlocking Agent* (sOA). The request can be submitted via a local interface, such as a hypervisor-specific shared memory implementation [68], [83], [95] or locally-terminated network endpoint [6], [34], [67]. The sOA predicts if there are enough resources to satisfy the request and, accordingly, grants or rejects the request. If the request is rejected, the local agent informs the global agent which then takes corrective actions (e.g., request scale-out or redistribute the load towards the overlocked VMs). In the background, each sOA monitors a server’s power and overlocking needs, and creates a profile to be periodically sent to the *Global Overlocking Agent* (gOA). The gOA uses the profiles to assign efficient per-server budgets. In turn, an sOA uses the assigned budget for admission control until the budget gets updated.

A. Workload-Aware Overlocking

Overview. SmartOClock extends the existing autoscaling interface with overlocking. A workload specifies the scale-up (start) and scale-down (stop) thresholds for overlocking. The overlocking hints can be inserted by developers after profiling or they can be automated using the existing tools for automatic instance scaling [11], [31], [64], [88], [99]. Like conventional autoscaling, the overlocking thresholds can be: metrics-based or schedule-based. Under *metrics-based* overlocking, workloads can use application metrics (e.g., tail latency, queue length) or resource utilization (e.g., CPU, network) to trigger overlocking. The granularity of application hints can be per-function in the case of tail latency or per-VM in the case of resource utilization. These metrics can then be monitored per- and across-VM instances for specified time intervals to meet an application’s goals. Additionally, workloads that have predictable times for high traffic (e.g., 9-10 AM local time) can use *schedule-based* thresholds. Finally, workloads can also use a combination of metrics- and schedule-based. Importantly, extending the autoscaling interface for overlocking enables using scaling out (creating new VMs) as a fallback mechanisms for when overlocking is not possible. The scale-out signal can also be triggered proactively by SmartOClock using predictions for the ability to overlock (see Section IV-D).

Adopting WI by cloud users. Although workload owners already carefully tune the metrics and thresholds for horizontal scaling, there is overhead in repeating the process for vertical scaling (overlocking). To ease adoption, SmartOClock can be extended to infer the overlocking thresholds. It can leverage workload historical data to determine scale-up values. The lifetime impact of overlocking can be factored in this analysis. For example, use P90 of historical value if overlocking

can be performed for 10% of the time only to comply with lifetime goals. The overclocking impact needs to be estimated to determine the scale-down value. An inaccurate estimate can either cause dithering if it is too close to the scale-up threshold or waste precious overclocking time if the estimate is too low. Performance models using low-level architectural counters can be used for the estimation. Workload owners can also leverage the inferred thresholds as an initial estimation.

B. Overclocking Admission Control

Overview. Naively granting overclocking requests (1) increases the chance of power capping events that deteriorate performance, and (2) accelerates wear out of server components. Instead, SmartOClock performs admission control for the overclocking requests based on *power and component lifetime impact predictions*. It predicts (1) the rack’s power draw to assess if overclocking will result in capping, and (2) the CPU utilization of VMs requesting overclocking to assess if overclocking them will exceed the lifetime budget. Using the predictions, SmartOClock decides (1) if the requested power and overclocking budgets can be reserved for a schedule-based workload, or (2) for how long a given VM can be overclocked under a metrics-based policy before needing corrective actions. Note that the power reservation is *soft*, the power can be taken by workloads outside of the system that do not need overclocking and SmartOClock needs to adjust.

Managing power. As observed in Section III, the power draw of racks and servers is highly predictable. Hence, the gOA and sOA continuously monitor the server and rack power draw and use the data gathered during monitoring to periodically (e.g., weekly) recompute the per-rack and per-server power templates. The templates are used to predict if the additional power of overclocking will trigger a capping event.

SmartOClock creates a power template using *per-day aggregation* of power draws across all weekdays in the prior week. The template represents a single day and the same template is used for predictions for all days in the following week. For example, the template’s value at 9AM is the median of rack’s power draw at 9AM across all five weekdays. A separate template is used for weekends. The intuition for this approach is that (1) using a coarse-grained measurement (e.g., the maximum over a week) is too conservative (i.e., it unnecessarily rejects many overclocking request) and (2) using fine-grained measurements (i.e., all power measurements from the prior week) is insufficiently robust to outliers (e.g., holidays during the prior week). Section V-B compares the accuracy of several template-creation strategies.

Managing lifetime impact from overclocking. A max time to overclock a component is obtained through an offline analysis with the vendors (e.g., 10% over a 5-year period). This analysis uses realistic, yet conservative, utilization of cloud components to determine the opportunity. The duration of individual overclockings can vary, but SmartOClock needs to honor the *total* overclocking time assumption to comply with component lifetime goals. This requirement is the same as for using turbo-boost on non-overclockable CPUs.

To get uniform overclocking over a component’s expected lifetime, SmartOClock divides the overall budget into epochs. The definition of an epoch is configurable (e.g., a day, week). Using a longer epoch, such as a week, enables assigning unused budgets from the weekend to the weekdays. Hence, SmartOClock defines an epoch to be a week and calculates per-weekday max overclocking time.

Each sOA ensures that the overclocked time-in-state of a component (e.g., per-core of a CPU) does not exceed the limit. Tracking and enforcement is per-server; an sOA uses mechanisms like Intel PMT [48] for the time-in-state tracking and denies overclocking requests if the budget is exhausted. Due to hardware heterogeneity, vendor-specific APIs are needed for the tracking; calling such APIs is already supported by operating systems (e.g., Intel PMT [49] and AMD HSMP [9] on Linux), and enforcement is via standard interfaces (e.g., CPPC [2] for CPU cores). For a predictable overclocking experience, an sOA also reserves overclocking budgets for scheduled requests. Unused budgets can be used by metrics-based overclocking and/or carried over to the next epoch.

C. Heterogeneous Power Budgets

Overview. SmartOClock splits the rack power budget *heterogeneously* amongst servers. Each sOA collect its server’s power draw and overclocking needs over time to create power and overclocking *templates*. The power template specifies a server’s draw at a given timestamp. The overclock template specifies the number of cores that *requested* and were *granted* overclocking. The sOAs periodically (e.g., weekly) exchange their templates with the gOA. The gOA combines power and overclocking templates of all sOAs and computes individual power budgets. It grants power credits to servers for periods when VMs are overclocked, per the reported template.

Power budget computation. The power budget computation happens in three phases. First, the gOA uses its power model to separate the server’s power into the regular and overclock power; the number of cores from the server’s overclocking template enable the gOA to discriminate the two portions. Second, the gOA assigns to each sOA the initial power budget that is equal to the server’s regular power draw. Finally, the gOA splits the remaining power headroom based on the overclocking requirements, i.e., servers with more overclocked cores in the past get larger extra power budgets for the future.

For example, a rack has two servers (X and Y) and 1.3kW power limit. Typical power draw without overclocking of X and Y at 9AM is 400W and 300W, respectively. Thus, the unused power is 600W. In addition, at 9AM, X and Y typically need to overclock 5 cores (extra 50W) and 10 cores (extra 100W), respectively. Based on this history, the gOA computes the power budgets for 9AM: for X $400W + \frac{50 \times 600}{50+100}W = 600W$, and for Y $300W + \frac{100 \times 600}{50+100}W = 700W$.

D. Decentralized Budget Enforcement

Overview. SmartOClock takes *decentralized* decisions by allowing servers to locally process overclocking requests from

their VMs. An sOA uses the server’s power profile to predict if overclocking will exceed the server’s power budget. As the budget computations rely on predictions, they may become stale. Thus, SmartOClock allows sOAs to explore beyond their initial assignments. Similarly, an sOA tracks the overclocking time of VMs and predicts if a VM will run out of budget. Then, to avoid missed SLOs, the sOA informs the global WI agent (via local) of the inability to overclock; in turn, the global WI agent can take corrective actions using the configured scale-out policies. Enabling local decisions is key for reactively handling activity bursts under metrics-based overclocking. The overclocking trigger by a WI agent is conveyed to the (local) sOA that can start/stop overclocking in order of a few milliseconds. Furthermore, if the assigned power budget is insufficient (*e.g.*, misprediction, due to change in load), then the sOA can independently explore a higher budget to maximize the extent (frequency) of overclocking.

Power budget enforcement. The gOA periodically sends the heterogeneously assigned power budgets to each sOA. Then, each sOA performs prioritized per-VM power management [57] via a feedback loop to control the server power draw while overclocking. For example, scheduled overclocking VMs can be of higher priority compared to unscheduled (metrics-based) ones. In the feedback loop, an sOA changes the frequency of the overclocked VMs per priority in discrete steps (*e.g.*, 100 MHz). Based on the impact of the last frequency change on the server’s power draw, the sOA either: (1) maintains the VMs at the current frequency (if $threshold \leq draw < limit$, where $threshold = limit - buffer$), (2) increases frequency by step size (if $draw < threshold$), or (3) reduces frequency by step size (if $draw > limit$). Prioritization enables overclocking the more important VMs to the maximum extent before less important VMs are overclocked.

Exploring beyond the local budgets. Due to mispredictions, the initial power allotment may become inefficient—some servers may consume less than predicted while others are limited by their budget and cannot overclock VMs to the maximum extent. Thus, SmartOClock allows sOAs to explore beyond their allocated power budgets. Specifically, on constrained servers, the sOA tries to gradually exceed the limit through two phases: *exploration* and *exploitation*.

Exploration. A sOA *conditionally* increases its budget by a step size (*e.g.*, 20W) that causes the feedback loop to start increasing the frequency of the overclocked VMs. If within a short timespan (*e.g.*, 30 seconds), the sOA does not receive any *warning* messages from the rack power capping system (run in the rack manager on each rack), then it further increases the budget. The sOA stops when all VMs are overclocked to the highest frequency or when it receives a warning message. The rack manager sends a warning message to all sOAs when the rack’s power draw reaches a *warning threshold* (*e.g.*, 95% of the rack’s power limit). An sOA ignores the message if it is not exploring. Otherwise, it reduces its budget by the step size and uses *exponential back-off* for the next exploration phase.

Exploitation. After establishing a safe power budget (*i.e.*, no

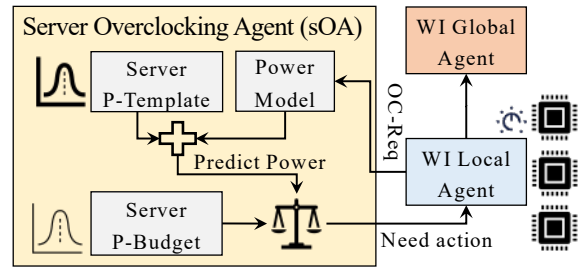


Fig. 11: Server Overclocking Agent in SmartOClock.

warning messages), a sOA enters the *exploitation* phase. In this phase, it uses the new budget to grant overclocking requests until either the *time to exploit* expires or upon receiving a *power capping event*. When the time to exploit expires, the sOA starts a new exploration phase if needed. Whereas, on a capping event, it goes back to its initial power budget.

Similarly, a sOA can explore beyond the local per-core overclocking budget. If a VM requires overclocking for longer than its assigned cores can sustain, the sOA will first overclock the VM’s assigned cores until their budget is exhausted. Then, the sOA explores rescheduling the VM on other cores in the server with available budget.

Managing resource exhaustion. When an overclocking request is rejected, the global WI agent takes corrective actions per an operator-chosen policy. A simple policy is to scale out while factoring the number of VMs that cannot be overclocked across a deployment (*e.g.*, create x new if y existing VMs cannot be overclocked). Figure 11 shows the operations performed by SmartOClock for managing power exhaustion. First, a sOA predicts when it will run out of power for overclocking. For this check, it first predicts the extra power from overclocking a given VM (for a worst-case CPU utilization). Next, via the template, it finds the time when the predicted extra power exceeds the server’s budget. It then sends a signal to the global WI agent if the time to exhaustion is within a configurable window (*e.g.*, 15 minutes). To minimize performance impact from a lack of overclocking, the length of the window should be greater than the time to scale out, so that overclocking is still available for the time it takes to scale out. Finally, this operation can be performed ahead of time for scheduled overclocking requests to protect workload SLOs. For metrics-based overclocking, the scale-up (overclocking) threshold can be set before scale-out, where SLOs would be missed if resources are inadequately provisioned after the scale-out threshold is exceeded. Setting an earlier scale-up threshold allows using overclocking to handle load spikes and enables reverting to scale-out if overclocking is not possible. Like power, an sOA also predicts the time to exhaustion of the overclocking budget and informs the global WI agent.

V. EVALUATION

To evaluate SmartOClock, we perform real-system experiments running cloud applications in an overclockable server cluster, and large-scale analysis using production traces.

A. Cluster-Level Experiments

Methodology. We implement SmartOClock and conduct the experiments on 36 overclockable servers (all 28 from one rack, and 8 from another during scale-out). Each server has a 64-core (128 threads) AMD EPYC 7763 CPU with customizations to facilitate overclocking experimentation. Its default max turbo frequency is 3.3GHz, which can be increased to 4.0 GHz on these custom parts for overclocking. The CPU is configured to operate in performance mode [82] and the active cores can steadily run at 4.0 GHz while TDP-unconstrained.

To set the load for each server, we take an example production rack from Azure. Based on the power traces of these production servers, we select which application to run in each individual server to mimic the *same* power profile. We run VMs hosting two open-source applications: (1) the latency-critical social network microservices (*SocialNet*) from DeathStarBench [32] and, (2) the throughput-optimized machine learning training (*MLTrain*) from FunctionBench [54]. In the power traces, 14 of the servers show constant high power while the other 14 show a diurnal pattern. For the first 14 servers, we use MLTrain and SocialNet for the other 14. The load for each benchmark instance is configured to mimic the power draw of the corresponding production server.

We define the per-server load in our experiments based on the production traces. As the profiled servers run different, independent, workloads, each server runs an independent set of SocialNet instances. Thus, there is no correlation in the power draw or loads across servers (*i.e.*, the load on one server does not affect the load on others). Auto-scaling is set for SocialNet based on its tail latency (initial count is 14). As in Section III, we set the SLO of each microservice to be 5 times its execution time on an unloaded system [26], [60], [73].

We compare SmartOClock with a *Baseline* system that does not scale horizontally (number of instances) nor vertically (core’s frequency), and *ScaleOut* and *ScaleUp* systems that only scale out/in and up/down, respectively, the number of SocialNet instances based on the observed tail. In the evaluation we use a metric-based overclocking policy, which is less predictable; experiments with a schedule-based policy show better results due to higher predictability.

Application performance. Figure 12 shows the P99 tail and average latency of SocialNet microservices in four environments. We group the 14 instances into three classes based on their load: Low, Medium, and High Load. Bars in the figure are the average across all instances with the same load level.

All systems perform equally well under low load. The impact on tail latency becomes prominent with increased load. Under high load, SmartOClock reduces the tail latency of Baseline, ScaleOut, and ScaleUp by 19.0%, 10.5%, and 8.9%.

The average latency of SmartOClock is lower than Baseline and ScaleUp, but slightly higher than ScaleOut. The reason is that, to reduce the application’s cost and prevent scaling out, SmartOClock operates for a longer time with higher latencies that are still below the SLO. However, SmartOClock significantly reduces the number of missed SLOs. The total

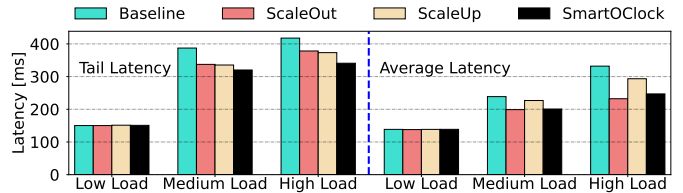


Fig. 12: P99 tail and average latencies of SocialNet services.

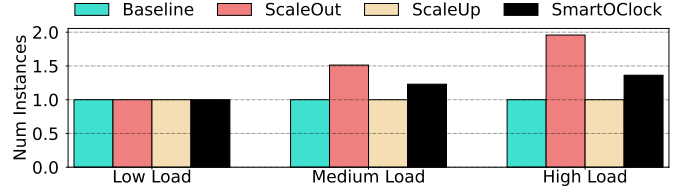


Fig. 13: Average number of SocialNet VMs varying load.

number of missed SLOs at high load is reduced by 26 \times , 4.8 \times , and 2.3 \times over Baseline, ScaleOut, and ScaleUp, respectively. These results show that overclocking (via ScaleUp or SmartOClock) reduces missed SLOs compared to ScaleOut. However, overclocking alone is insufficient at higher loads as evidenced by the greater missed SLOs with ScaleUp, despite it overclocking for 5x longer. A combination of ScaleUp and ScaleOut via SmartOClock provides the best performance. Finally, SmartOClock reacts fast to sudden workload shifts and keeps the application performance within its SLO: even on servers that triggered overclocking for more than 140 times within 5 minutes, SmartOClock did not miss any deadlines.

Cost. Performance improvements from SmartOClock result in cost savings for the users as they need to pay for fewer VMs. Figure 13 shows the average number of concurrently active VM instances for each environment over the entire run. Under high load, SmartOClock saves substantial cost by reducing the number of required instances by 30.4% over ScaleOut.

Energy consumption. Figure 14 shows normalized (1) per-single-server energy consumption under low, medium, and high load, and (2) total energy consumption of the system. Note that ScaleOut and SmartOClock are the only systems that meet SLOs. As the load increases, SmartOClock frequently overclocks cores, which increases the per-server energy consumption. However, as it uses fewer instances, the total energy consumption is reduced by 10% on average over ScaleOut. The savings are larger if we only consider servers running latency-critical microservices — 23% on average over ScaleOut.

Power-constrained environments. We evaluate SmartOClock’s overclocking admission control and heterogeneous power budgeting under constraints. We reduce the rack’s limit and measure the performance in two systems: NaiveOClock and SmartOClock. NaiveOClock grants all overclocking requests and on a power capping event splits the rack’s budget equally among the servers. SmartOClock reduces the SocialNet tail latency by 6.7% and 8.4% for medium and high loads, respectively, and improves the MLTrain throughput by 10.4%.

Overclocking-constrained environments. To evaluate SmartOClock’s proactive scale-out, we restrict the overclocking

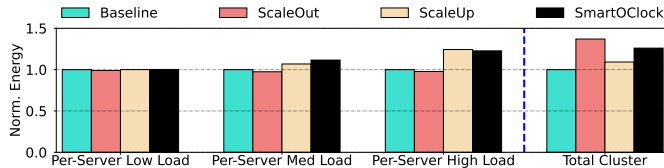


Fig. 14: Normalized per-single-server energy.

budget and measure the number of missed SLOs with and without proactive scaling. As we reduce the budget to 75%, 50%, and 25% of its initial value, reactive scale-out misses the SLO for 5.0%, 6.1%, and 7.2% of time, while SmartOClock’s proactive approach eliminates all SLO violations.

B. Large-Scale Simulations

Methodology. We use production traces of dedicated racks running Microsoft’s productivity and collaboration services (see Section III) from multiple datacenters. Each datacenter deployment is composed of hundreds of racks and a few thousand servers with either Intel or AMD CPUs. Each workload’s VMs are spread across servers and racks. The traces include rack and server power, and VM-level CPU utilization. All data is collected for 6 weeks (April 10th - May 12th, 2023), at a 5-minute granularity. Overclocking requirements (e.g., time of day) are obtained from the workload operators.

We develop a discrete event simulator to evaluate SmartOClock. Models are used to estimate the power impact of overclocking; CPU utilization and core frequency are the input. We validate the model for each server generation.

We compare SmartOClock to (1) *Central* – an oracle with a global view of power draw that can precisely decide if an overclocking request will result in capping, (2) *NaiveOClock* – a system that grants all overclocking requests, (3) *NoFeedback* – a system that adheres to the per-server power budgets with no exploration beyond, and (4) *NoWarning* – a system that allows exploring but with no warnings. The servers go back to their initial power budget on a capping event.

Overclocking success and power capping. Table I shows the results: (1) number of power capping events in each system normalized to Central, (2) percentage of successful overclocking requests, (3) performance penalty of capping on non-overclocked VMs, and (4) normalized performance over Baseline. We define the performance penalty and improvement as reduction and increase in VM frequency compared to the Baseline (max turbo), respectively. Clusters are split into three groups based on power draw: *High*, *Medium*, and *Low-Power*.

First, naively granting overclocking requests causes many power capping events. NaiveOClock causes 118.6 \times , 36.6 \times , and 14.0 \times more events than Central in High, Medium, and Low-Power clusters, respectively. In contrast, SmartOClock lowers the events by 18.9 \times in High-Power clusters via prediction for admission control. Adding the warning messages efficiently controls overclocking beyond a server’s budget: it reduces the number of events over NoWarning by up to 4.3 \times .

Second, SmartOClock successfully grants majority of overclocking requests. It is within 4%, 3%, and 1% of the success

TABLE I: Comparison of SmartOClock to different baselines.

System	Norm. # of Power Caps	Successful OClock Reqs	Penalty on Power Cap	Norm. Performance
High-Power Clusters				
Central	1.0	92%	21%	1.186
NaiveOClock	118.6	55%	34%	0.963
NoFeedback	5.5	72%	22%	1.122
NoWarning	27.4	81%	23%	1.081
SmartOClock	6.3	89%	22%	1.164
Medium-Power Clusters				
Central	1.0	96%	11%	1.195
NaiveOClock	36.6	79%	19%	1.022
NoFeedback	3.4	83%	11%	1.163
NoWarning	7.2	87%	12%	1.160
SmartOClock	3.9	93%	11%	1.185
Low-Power Clusters				
Central	1.0	99%	1%	1.208
NaiveOClock	14.0	99%	5%	1.172
NoFeedback	1.0	98%	1%	1.205
NoWarning	1.1	99%	2%	1.205
SmartOClock	1.0	99%	1%	1.208

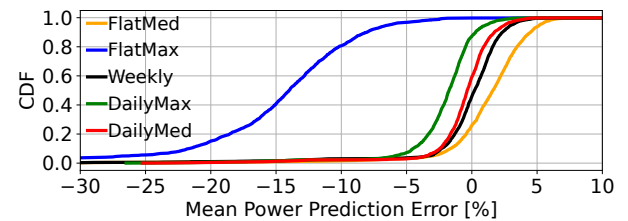


Fig. 15: CDF of mean power prediction for each technique.

rate of an oracle Central system in High, Medium, and Low-Power clusters, respectively. The feedback-loop for exploring beyond the per-server budget is important: SmartOClock has up to 1.24 \times higher success rate than NoFeedback approach.

Finally, heterogeneous power distribution by SmartOClock reduces the performance penalty from power capping events. All systems bar NaiveOClock employ this optimization. The heterogeneous power budgets reduce the performance penalty due to power capping events over NaiveOClock by 1.62 \times and 1.72 \times in High and Medium-Power clusters, respectively.

Power predictions accuracy. Figure 15 shows the CDF of prediction accuracy for computing the power templates. *FlatMed* and *FlatMax* use a constant prediction: a median or maximum of all prior measurements. FlatMed is opportunistic and underpredicts power, leading to high P99 prediction errors. Whereas, FlatMax is conservative and overpredicts power, resulting in negative prediction errors at low percentiles.

Weekly uses a time series of power measurements from the previous week for predictions in the following week. It is impacted by outliers since it treats each day separately: a few hours may behave differently due to the unexpected events. Thus, at high percentiles, its prediction error can be significant.

Finally, *DailyMed* and *DailyMax*, aggregate the power measurements across a week to represent a single, typical, day. The templates are time series of median or maximum values. DailyMed, used in SmartOClock, has the highest accuracy.

C. Experiments with Production Services

We evaluate overclocking *Service B* and *C* under production load. Each service consumes hundreds of virtual cores across

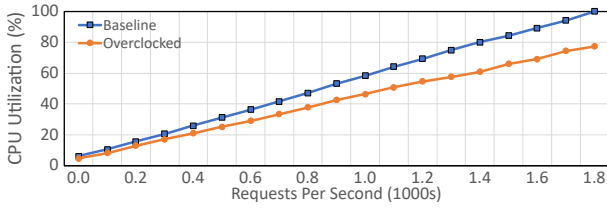


Fig. 16: Impact of overclocking Service B.

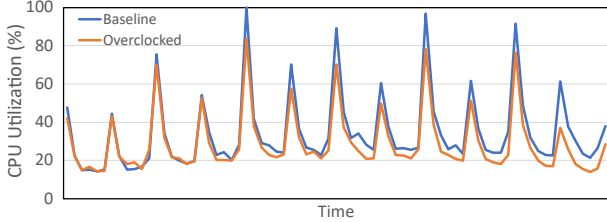


Fig. 17: Impact of overclocking Service C.

tens of VMs. The deployment resource usage is similar to Figure 1 and the SLOs are consistent with each service’s goals.

Figure 16 shows the average CPU utilization of *Service B*’s VMs for different request rates (bucketized by 0.1 due to live load variability). Overclocking reduces CPU utilization of VMs by 23% at a peak of 1.8k requests per second (RPS); the baseline operates at turbo (3.3 GHz). Alternatively, for the same CPU utilization, baseline can service 1.4k RPS vs 1.8k (28% higher) with overclocking. Figure 17 shows that overclocking reduces *Service C*’s 5-minute peaks over a weekday by 16%. The deployment load is similar on both days. Both results show the opportunity to down-provision while meeting the performance SLOs. Finally, overclocking enables servicing 25% additional load by *Service A* VMs under synthetic traffic; production experiments are being setup.

VI. LESSONS FROM PRODUCTION DEPLOYMENT

We built a first-of-its-kind 2-rack (56 servers) overclockable cluster at a cloud provider for CPU overclocking in production. Our deployment does not yet include cluster-wide coordination. Here we present lessons from the deployment.

Motivation for building a cluster. Although CPU overclocking can provide substantial performance and cost benefits, a comprehensive analysis (*e.g.*, TCO reduction, revenue increase) is needed for introducing hardware features at scale. Projecting improvements is challenging due to workload-specific variations, as previous work shows [51]. Furthermore, evaluating in a lab environment is not possible for even Microsoft’s internal workloads due to software dependencies (*e.g.*, deployment framework) and security concerns that prevent experimentation with production traffic. Thus, building an overclockable cluster was imperative.

Using experimental hardware in production datacenters. Azure has a rigid process for ensuring stability (*e.g.*, thermal limitations), reliability (*e.g.*, firmware errors), and high performance for hardware deployment at scale that adds overhead for limited-scale experimentation. To address, we retrofitted overclocking onto existing hardware by installing

overclockable CPUs and firmware updates (*e.g.*, BIOS) on already-deployed servers in a production datacenter. We also bypassed software checks that remove servers with unexpected configuration. A drawback of our approach is that the platform (*e.g.*, motherboard, cooling) is not optimized for overclocking, leading to thermal and max current throttling under heavy loads that affect performance. Additionally, we provisioned adequate power for the racks to avoid capping; the limits are lowered for power management evaluations.

Experiments with first-party workloads. Since the cluster contains experimental hardware, we enforce strict admission control. However, this policy led to atypical workload placements that impact overclocking. Typically servers house VMs from various workloads due to dynamic cloud environments and scheduler efforts to optimize resource utilization [40], [90], but our policy caused VMs from the same workload to occupy entire servers. Although this impacts overclocking efficiency, it is useful for conservative benefit estimation.

Finer-grained overclocking. SmartOClock can overclock individual VMs but first-party operators want finer-grained overclocking (*e.g.*, containers in VMs). Although overclocking VMs still works, it is inefficient because of the higher power and reliability impact. Since containers are scheduled inside a guest VM without host visibility, we need guest participation for finer-grained overclocking. However, unsupervised control of frequency by guests can compromise reliability and power management. We are exploring a safe and efficient solution.

Hardware support for overclocking. Overclocking lifetime budgets can be improved with *wear-out counters* that indicate how a component’s (*e.g.*, CPU core) lifetime is impacted by utilization (voltage) and operating temperatures. SmartOClock can use wearout counters to upgrade from a conservative offline model to a *per-part* online calculation for safety.

Furthermore, the prioritized feedback loop for managing power while overclocking can be offloaded to the hardware for efficiency. We are extending the ACPI [2] CPPC interface to configure VM priority while scheduling (no affinity) on CPU cores. The firmware can use these priorities to assign per-core performance (frequency) while managing power.

Vendor engagements to enable overclocking. As overclocking is enabled by under-utilization (Section III), instead of overclocking, vendors (*e.g.*, Intel, AMD) inquire about designing a CPU with revised time-in-state assumptions for offline certification. However, this is still inefficient as it does not leverage the utilization variability from workload demands (with and without overclocking) and temperature fluctuations on ageing at cloud scale. Using wear-out counters to track the usage impact on ageing does not have these limitations.

Furthermore, we are working with the vendors to ensure all cores can hit a minimum-desired overclocking frequency (*e.g.*, 15-20% beyond max turbo). Some cores can run faster, but this variability is not exposed on server CPUs (even for turbo); we are exploring bringing mechanisms from client CPUs (*e.g.*, ACPI CPPC preferred cores [2]) to leverage this variability.

Overclocking beyond CPUs. SmartOClock is a general framework and its principles can be easily applied for overclocking any server component. Our initial focus was CPU since it provides the highest benefits, but we have started exploring overclocking of other components (*e.g.*, GPU).

Silent data corruption (SDC). Prior work shows the risks from SDC at scale [28], [41], [92]. Although overclocking can aggravate error rates due to aggressive circuit timing and sudden voltage drops, our extensive lab and production experiments do not show an increase in errors, with frequencies $\sim 20\%$ beyond max turbo; this is inline with a prior work [51]. Nonetheless, for safety, we work with the vendors to define a max overclocking frequency. Furthermore, techniques from the SDC work can also be used for added safeguarding.

VII. RELATED WORK

Computational sprinting. Extensive research [15], [16], [18], [30], [55], [63], [74], [78], [79], [100], [101] has explored computational sprinting (*i.e.*, boosting CPU frequency for short periods). Mechanisms like game theory [30], formal control [77], and performance modeling [74] have been proposed to manage sprinting. Researchers have also investigated efficiency factors like resource interference [63], power availability [16], processor design [38], and cooling [51]. However, none of these works holistically address the overclocking challenges in the cloud. They either focus on a single-server setup, assume a transparent-box knowledge of the applications, or overlook multi-tenancy on a server or rack.

The closest related work is Computational Sprinting Game (CSG) [30]. There are two major differences between CSG and SmartOClock. First, CSG leverages turbo and is constrained by thermal/power limits. In contrast, overclocking also affects reliability whose time scales are orders of magnitude (months/years) more than for power/thermal (minutes). It is nontrivial to add reliability under CSG when evaluating sprint utility. In contrast, SmartOClock uses epochs to divide the overclocking budget across coarse-grain time scales (days) that local agents enforce. Second, a lack of sprinting/overclocking (of even a few VMs) can impact the SLOs of workloads that underprovision while relying on sprinting to handle their peaks. Therefore, a mitigation mechanism to protect performance is needed when sprinting is unavailable, a problem not addressed by CSG. Section V presents the impact of proactive scaleout by SmartOClock to protect workload SLOs.

Undervolting. Prior work has proposed decreasing the voltage for a frequency below its safe marginal value for reducing power [12], [13], [17], [29], [52], [75]. However, undervolting can introduce instability and pipeline (*i.e.*, timing) errors, thereby necessitating hardware designers to add mechanisms for fault tolerance. For example, Razor [29] uses additional latches that run on a delayed clock in vulnerable paths to detect/recover from errors. This body of work is complementary and can create additional power and component lifetime headroom (reduced wearout from lower voltage) for overclocking.

Datacenter power management. Prior work has proposed oversubscription through leveraging statistical properties of concurrent power usage across servers [14], [37], [42], [56], [57], [59], [80], [84], [94] to improve datacenter power utilization and save costs. These works are complementary and influence our non-overclocked baseline. Azure leverages policies based on these prior works to oversubscribe power. The policies factor the power demand from turbo for meeting the performance SLAs [7], [36], [72] and prioritized throttling [57], [59] is used to protect (turbo) performance of critical workloads under rare power capping events.

Naively adding overclocking to the baseline power utilization increases the probability of power capping events. Increasing provisioned power cannot be used to address this problem due to the TCO impact, especially when turbo is sufficient to meet a provider’s performance SLAs. Consequently, an overclocking system can only leverage unutilized power while meeting workload SLOs when overclocking is not possible; problems not addressed by the prior power management work.

Workload intelligence. Research has leveraged workload awareness to optimize performance, energy consumption, and cost [19]–[21], [44], [53], [87], [99], [102]. Sinan [99] uses ML models to allocate resources per microservice tier for minimizing cost while maintaining latency targets. Re-Tail [19], Rubik [53], Adrenaline [43], and Gemini [102] use application-specific features to predict optimal per-request frequencies, reducing power draw while meeting SLOs. Resource Central [21] gathers VM telemetry, learns VM behaviors offline, and provides online predictions for various resource managers. We propose a clean interface for cloud workloads to provide the necessary signals for overclocking without compromising their opaque-box implementations.

VIII. CONCLUSION

In this paper, we proposed SmartOClock, the first distributed overclocking management platform for cloud environments. SmartOClock enables cloud providers to offer overclocking to workloads through four novel features: workload intelligence, prediction-based admission control, heterogeneous power budgeting, and decentralized enforcement. Our evaluation shows that SmartOClock reduces the tail latency by 8.9% and the application cost by 30.4%. We also discussed lessons from building an overclockable cluster in Azure. We conclude that carefully-managed overclocking has enormous potential to improve workload performance while saving cost.

IX. ACKNOWLEDGEMENTS

We would like to thank the reviewers for helping us improve this paper. We also thank Saurabh Agrawal and Georgia Modoran for help with CPU reliability analysis, Serena Zhao and Juan Perez with the cluster buildout, Jim Kleewein, Nilesh Oswal, Sergey Anikin and Rui Liang for helping run production experiments, and Paul Artman and Anil Harwani from AMD for enabling us with overclockable CPUs.

REFERENCES

- [1] S. I. Abrita, M. Sarker, F. Abrar, and M. A. Adnan, "Benchmarking VM Startup Time in the Cloud," in *Benchmarking, Measuring, and Optimizing: First BenchCouncil International Symposium*, 2019.
- [2] ACPI Specification Revision Committee, "Advanced configuration and power interface specification," 2022. [Online]. Available: <https://uefi.org/specifications>
- [3] P. Alcorn, "How to Overclock Your CPU: CPU Overclocking Impact on Lifespan and Reliability," May 2023. [Online]. Available: <https://www.tomshardware.com/how-to/how-to-overclock-a-cpu#section-cpu-overclocking-impact-on-lifespan-and-reliability>
- [4] Amazon AWS, "AWS Auto Scaling," April 2024. [Online]. Available: <https://aws.amazon.com/autoscaling/>
- [5] Amazon AWS, "AWS Lambda," April 2024. [Online]. Available: <https://aws.amazon.com/lambda/>
- [6] Amazon AWS, "Instance metadata and user data," April 2024. [Online]. Available: <https://docs.aws.amazon.com/AWSEC2/latest/UserGuide/ec2-instance-metadata.html>
- [7] Amazon Web Services, "Processor state control for your EC2 instance," April 2024. [Online]. Available: https://docs.aws.amazon.com/AWSEC2/latest/UserGuide/processor_state_control.html
- [8] AMD, "CPU warranty terms," April 2024. [Online]. Available: <https://www.amd.com/system/files/documents/processor-warranty-update.pdf>
- [9] AMD, "Host System Management Port (HSMP)," April 2024. [Online]. Available: https://github.com/amd/amd_hsmp
- [10] AMD, "Turbo Core Technology," April 2024. [Online]. Available: <https://www.amd.com/en/technologies/turbo-core>
- [11] A. F. Baarzi and G. Kesidis, "SHOWAR: Right-Sizing And Efficient Scheduling of Microservices," in *Proceedings of the 12th Symposium on Cloud Computing (SoCC '21)*, 2021.
- [12] A. Bacha and R. Teodorescu, "Dynamic reduction of voltage margins by leveraging on-chip ECC in Itanium II processors," in *Proceedings of the 40th Annual International Symposium on Computer Architecture (ISCA '13)*, 2013.
- [13] R. Bertran, A. Buyuktosunoglu, P. Bose, T. J. Slegel, G. Salem, S. Carey, R. F. Rizzolo, and T. Strach, "Voltage Noise in Multi-Core Processors: Empirical Characterization and Optimization Opportunities," in *Proceedings of the 47th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO '14)*, 2014.
- [14] A. Bhattacharya, D. Culler, A. Kansal, S. Govindan, and S. Sankar, "The need for speed and stability in data center power capping," in *Proceedings of the International Green Computing Conference (IGCC '12)*, 2012.
- [15] H. Cai, Q. Cao, F. Sheng, Y. Yang, C. Xie, and L. Xiao, "ESprint: QoS-Aware Management for Effective Computational Sprinting in Data Centers," in *Proceedings of the 19th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (CCGRID '19)*, 2019.
- [16] H. Cai, X. Zhou, Q. Cao, H. Jiang, F. Sheng, X. Qi, J. Yao, C. Xie, L. Xiao, and L. Gu, "GreenSprint: Effective Computational Sprinting in Green Data Centers," in *Proceedings of the IEEE International Parallel and Distributed Processing Symposium (IPDPS '18)*, 2018.
- [17] K. K. Chang, A. G. Yağlıkçı, S. Ghose, A. Agrawal, N. Chatterjee, A. Kashyap, D. Lee, M. O'Connor, H. Hassan, and O. Mutlu, "Understanding Reduced-Voltage Operation in Modern DRAM Devices: Experimental Characterization, Analysis, and Mechanisms," *Proceedings of the ACM on Measurements and Analysis of Computer Systems*, 2017.
- [18] J. Charles, P. Jassi, N. S. Ananth, A. Sadat, and A. Fedorova, "Evaluation of the Intel® Core™ i7 Turbo Boost feature," in *Proceedings of the IEEE International Symposium on Workload Characterization (IISWC '09)*, 2009.
- [19] S. Chen, A. Jin, C. Delimitrou, and J. F. Martínez, "ReTail: Opting for Learning Simplicity to Enable QoS-Aware Power Management in the Cloud," in *Proceedings of the IEEE 28th International Symposium on High-Performance Computer Architecture (HPCA '22)*, 2022.
- [20] Z. Chen, J. Hu, G. Min, A. Y. Zomaya, and T. El-Ghazawi, "Towards Accurate Prediction for High-Dimensional and Highly-Variable Cloud Workloads with Deep Learning," *IEEE Transactions on Parallel and Distributed Systems*, vol. 31, no. 4, 2020.
- [21] E. Cortez, A. Bonde, A. Muzio, M. Russinovich, M. Fontoura, and R. Bianchini, "Resource Central: Understanding and Predicting Workloads for Improved Resource Management in Large Cloud Platforms," in *Proceedings of the 26th Symposium on Operating Systems Principles (SOSP '17)*, 2017.
- [22] H. David, E. Gorbato, U. R. Hanebutte, R. Khanna, and C. Le, "RAPL: Memory power estimation and capping," in *Proceedings of the ACM/IEEE International Symposium on Low-Power Electronics and Design (ISLPED '10)*, 2010.
- [23] W. R. Davis, C. Shaw, and A. R. Hassan, "How to write a compact reliability model with the open model interface (omi)," in *2020 IEEE International Reliability Physics Symposium (IRPS)*, 2020, pp. 1–2.
- [24] J. Dean and L. A. Barroso, "The Tail at Scale," *Communications of the ACM*, vol. 56, pp. 74–80, 2013.
- [25] C. Delimitrou and C. Kozyrakis, "Quasar: Resource-Efficient and QoS-Aware Cluster Management," in *Proceedings of the 19th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS '14)*, 2014.
- [26] C. Delimitrou and C. Kozyrakis, "Amdahl's Law for Tail Latency," *Commun. ACM*, vol. 61, no. 8, jul 2018.
- [27] D. DiMaria and J. Stathis, "Non-arrhenius temperature dependence of reliability in ultrathin silicon dioxide films," *Applied Physics Letters*, 1999.
- [28] H. D. Dixit, S. Pendharkar, M. Beadon, C. Mason, T. Chakravarthy, B. Muthiah, and S. Sankar, "Silent Data Corruptions at Scale," *CoRR*, vol. abs/2102.11245, 2021. [Online]. Available: <https://arxiv.org/abs/2102.11245>
- [29] D. Ernst, N. S. Kim, S. Das, S. Pant, R. Rao, T. Pham, C. Ziesler, D. Blaauw, T. Austin, K. Flautner, and T. Mudge, "Razor: a low-power pipeline based on circuit-level timing speculation," in *Proceedings. 36th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO '03)*, 2003.
- [30] S. Fan, S. M. Zahedi, and B. C. Lee, "The Computational Sprinting Game," in *Proceedings of the 21st International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS '16)*, 2016.
- [31] Y. Gan, M. Liang, S. Dev, D. Lo, and C. Delimitrou, "Sage: practical and scalable ML-driven performance debugging in microservices," in *Proceedings of the 26th International Conference on Architectural Support for Programming Languages and Operating Systems*, 2021.
- [32] Y. Gan, Y. Zhang, D. Cheng, A. Shetty, P. Rathi, N. Katarki, A. Bruno, J. Hu, B. Ritchken, B. Jackson, K. Hu, M. Pancholi, Y. He, B. Clancy, C. Colen, F. Wen, C. Leung, S. Wang, L. Zaruvisky, M. Espinosa, R. Lin, Z. Liu, J. Padilla, and C. Delimitrou, "An Open-Source Benchmark Suite for Microservices and Their Hardware-Software Implications for Cloud & Edge Systems," in *Proceedings of the 24th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS '19)*, 2019.
- [33] Google Cloud, "Autoscaling groups of instances," 2023. [Online]. Available: <https://cloud.google.com/compute/docs/autoscaler/>
- [34] Google Cloud, "About VM metadata," April 2024. [Online]. Available: <https://cloud.google.com/compute/docs/metadata/overview>
- [35] Google Cloud, "Google Cloud Functions," April 2024. [Online]. Available: <https://cloud.google.com/functions>
- [36] Google Compute Platform, "CPU platforms," April 2024. [Online]. Available: <https://cloud.google.com/compute/docs/cpu-platforms>
- [37] S. Govindan, J. Choi, B. Urgaonkar, A. Sivasubramaniam, and A. Baldini, "Statistical profiling-based techniques for effective power provisioning in data centers," in *Proceedings of the 4th ACM European Conference on Computer Systems (EuroSys '09)*, 2009.
- [38] B. Greskamp and J. Torrellas, "Paceline: Improving Single-Thread Performance in Nanoscale CMPs through Core Overclocking," in *Proceedings of the 16th International Conference on Parallel Architecture and Compilation Techniques (PACT '07)*, 2007.
- [39] J. Guo, Z. Chang, S. Wang, H. Ding, Y. Feng, L. Mao, and Y. Bao, "Who Limits the Resource Efficiency of My Datacenter: An Analysis of Alibaba Datacenter Traces," in *Proceedings of the IEEE/ACM 27th International Symposium on Quality of Service (IWQoS '19)*, 2019.
- [40] O. Hadary, L. Marshall, I. Menache, A. Pan, E. E. Greeff, D. Dion, S. Dorminey, S. Joshi, Y. Chen, M. Russinovich, and T. Moscibroda, "Protean: VM Allocation Service at Scale," in *Proceedings of the 14th USENIX Symposium on Operating Systems Design and Implementation (OSDI '20)*, 2020.
- [41] P. H. Hochschild, P. J. Turner, J. C. Mogul, R. K. Govindaraju, P. Ranganathan, D. E. Culler, and A. Vahdat, "Cores that don't count," in *Proceedings of the Workshop on Hot Topics in Operating Systems (HotOS '21)*, 2021.

- [42] C.-H. Hsu, Q. Deng, J. Mars, and L. Tang, "SmoothOperator: Reducing Power Fragmentation and Improving Power Utilization in Large-Scale Datacenters," in *Proceedings of the 23rd International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS '18)*, 2018.
- [43] C.-H. Hsu, Y. Zhang, M. A. Laurenzano, D. Meisner, T. Wenisch, J. Mars, L. Tang, and R. G. Dreslinski, "Adrenaline: Pinpointing and reining in tail queries with quick voltage boosting," in *Proceedings of the IEEE 21st International Symposium on High Performance Computer Architecture (HPCA '15)*, 2015.
- [44] Q. Hu, P. Sun, S. Yan, Y. Wen, and T. Zhang, "Characterization and Prediction of Deep Learning Workloads in Large-Scale GPU Datacenters," in *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis (SC '21)*, 2021.
- [45] IBM Cloud, "IBM Cloud Functions," April 2024. [Online]. Available: <https://cloud.ibm.com/functions/>
- [46] IBM Cloud, "'Scaling stateful and stateless services'," April 2024. [Online]. Available: <https://www.ibm.com/docs/en/cloud-app-management/2019.3.0?topic=sizing-scaling-stateless-stateful-services>
- [47] Intel, "CPU warranty terms," April 2024. [Online]. Available: https://www.intel.com/content/dam/support/us/en/documents/processors/Limited_Warranty_8.5x11_for_Web_English.pdf
- [48] Intel, "Intel Platform Analysis Technology," April 2024. [Online]. Available: <https://www.intel.com/content/www/us/en/developer/topic-technology/platform-analysis-technology/overview.html>
- [49] Intel, "Platform Monitoring Technology Telemetry (PMT)," April 2024. [Online]. Available: <https://github.com/intel/Intel-PMT>
- [50] Intel, "What Is Intel® Turbo Boost Technology?" April 2024. [Online]. Available: <https://www.intel.com/content/www/us/en/gaming/resources/turbo-boost.html>
- [51] M. Jalili, I. Manousakis, I. Goiri, P. A. Misra, A. Raniwala, H. Alissa, B. Ramakrishnan, P. Tuma, C. Belady, M. Fountoura, and R. Bianchini, "Cost-Efficient Overclocking in Immersion-Cooled Datacenters," in *Proceedings of the 48th Annual International Symposium on Computer Architecture (ISCA '21)*, 2021.
- [52] M. Kaliorakis, A. Chatzidimitriou, G. Papadimitriou, and D. Gizopoulos, "Statistical Analysis of Multicore CPUs Operation in Scaled Voltage Conditions," *IEEE Computer Architecture Letters*, 2018.
- [53] H. Kasture, D. B. Bartolini, N. Beckmann, and D. Sanchez, "Rubik: Fast analytical power management for latency-critical systems," in *Proceedings of the 48th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO '15)*, 2015.
- [54] J. Kim and K. Lee, "FunctionBench: A Suite of Workloads for Serverless Cloud Function Service," in *Proceedings of the IEEE 12th International Conference on Cloud Computing (CLOUD '19)*, 2019.
- [55] S. Kondguli and M. Huang, "A Case for a More Effective, Power-Efficient Turbo Boosting," *ACM Transactions on Architecture and Code Optimization (TACO '18)*, vol. 15, no. 1, 2018.
- [56] V. Kontorinis, L. E. Zhang, B. Aksanli, J. Sampson, H. Homayoun, E. Pettis, D. M. Tullsen, and T. S. Rosing, "Managing Distributed Ups Energy for Effective Power Capping in Data Centers," in *Proceedings of the 39th Annual International Symposium on Computer Architecture (ISCA '12)*, 2012.
- [57] A. G. Kumbhare, R. Azimi, I. Manousakis, A. Bonde, F. Frujeri, N. Mahalingam, P. A. Misra, S. A. Javadi, B. Schroeder, M. Fountoura, and R. Bianchini, "Prediction-Based Power Oversubscription in Cloud Platforms," in *Proceedings of the USENIX Annual Technical Conference (USENIX ATC '21)*, 2021.
- [58] Y.-H. Lee, N. R. Mielke, W. McMahon, Y.-L. R. Lu, and S. Pae, "Thin-gate-oxide breakdown and cpu failure-rate estimation," *IEEE Transactions on Device and Materials Reliability*, vol. 7, no. 1, pp. 74–83, 2007.
- [59] S. Li, X. Wang, X. Zhang, V. Kontorinis, S. Kodakara, D. Lo, and P. Ranganathan, "Thunderbolt: Throughput-Optimized, Quality-of-Service-Aware Power Capping at Scale," in *Proceedings of the 14th USENIX Symposium on Operating Systems Design and Implementation (OSDI '20)*, 2020.
- [60] Z. Li, L. Zheng, Y. Zhong, V. Liu, Y. Sheng, X. Jin, Y. Huang, Z. Chen, H. Zhang, J. E. Gonzalez, and I. Stoica, "AlpaServe: Statistical Multiplexing with Model Parallelism for Deep Learning Serving," in *Proceedings of the 17th USENIX Symposium on Operating Systems Design and Implementation (OSDI '23)*, 2023.
- [61] H. Lim, D. Han, D. G. Andersen, and M. Kaminsky, "MICA: A Holistic Approach to Fast In-Memory Key-Value Storage," in *Proceedings of the 11th USENIX Symposium on Networked Systems Design and Implementation (NSDI '14)*, 2014.
- [62] Q. Liu and Z. Yu, "The Elasticity and Plasticity in Semi-Containerized Co-Locating Cloud Workload: A View from Alibaba Trace," in *Proceedings of the 9th Symposium on Cloud Computing (SoCC '18)*, 2018.
- [63] D. Lo and C. Kozyrakis, "Dynamic management of TurboMode in modern multi-core chips," in *Proceedings of the IEEE 20th International Symposium on High Performance Computer Architecture (HPCA '14)*, 2014.
- [64] S. Luo, H. Xu, K. Ye, G. Xu, L. Zhang, G. Yang, and C. Xu, "The power of prediction: microservice auto scaling via workload learning," in *Proceedings of the 13th Symposium on Cloud Computing (SoCC '22)*, 2022.
- [65] J. Lyu, M. You, C. Irvine, M. Jung, T. Narmore, J. Shapiro, L. Marshall, S. Samal, I. Manousakis, L. Hsu, P. Subbarayalu, A. Raniwala, B. Warrior, R. Bianchini, B. Schroeder, and D. S. Berger, "HyraX: Fail-in-Place Server Operation in Cloud Platforms," in *Proceedings of the 17th USENIX Symposium on Operating Systems Design and Implementation (OSDI '23)*, 2023.
- [66] D. Marcon, T. Kauerauf, F. Medjdoub, J. Das, M. Van Hove, P. Srivastava, K. Cheng, M. Leys, R. Mertens, S. Decoutere, G. Meneghesso, E. Zanoni, and G. Borghs, "A comprehensive reliability investigation of the voltage-, temperature- and device geometry-dependence of the gate degradation on state-of-the-art GaN-on-Si HEMTs," in *Proceedings of the 2010 International Electron Devices Meeting*, 2010.
- [67] Microsoft Azure, "Azure Instance Metadata Service," April 2024. [Online]. Available: <https://learn.microsoft.com/en-us/azure/virtual-machines/instance-metadata-service>
- [68] Microsoft Azure, "Data Exchange: Using key-value pairs to share information between the host and guest on Hyper-V," April 2024. [Online]. Available: [https://learn.microsoft.com/en-us/previous-versions/windows/it-pro/windows-server-2012-R2-and-2012/dn798287\(v=ws.11\)](https://learn.microsoft.com/en-us/previous-versions/windows/it-pro/windows-server-2012-R2-and-2012/dn798287(v=ws.11))
- [69] Microsoft Azure, "Introduction to Auto Scaling," April 2024. [Online]. Available: <https://learn.microsoft.com/en-us/azure/service-fabric/service-fabric-cluster-resource-manager-autoscaling>
- [70] Microsoft Azure, "Microsoft Azure Functions," April 2024. [Online]. Available: <https://azure.microsoft.com/en-gb/services/functions/>
- [71] Microsoft Azure, "Overview of autoscale in Azure," April 2024. [Online]. Available: <https://learn.microsoft.com/en-us/azure/azure-monitor/autoscale/autoscale-overview>
- [72] Microsoft Azure, "Virtual Machine series," April 2024. [Online]. Available: <https://azure.microsoft.com/en-us/pricing/details/virtual-machines/series>
- [73] A. Mirhosseini and T. Wenisch, "μSteal: A Theory-Backed Framework for Preemptive Work and Resource Stealing in Mixed-Criticality Microservices," in *Proceedings of the ACM International Conference on Supercomputing (ICS '21)*, 2021.
- [74] N. Morris, C. Stewart, L. Chen, R. Birke, and J. Kelley, "Model-Driven Computational Sprinting," in *Proceedings of the 13th ACM European Conference on Computer Systems (EuroSys '18)*, 2018.
- [75] G. Papadimitriou, M. Kaliorakis, A. Chatzidimitriou, D. Gizopoulos, P. Lawthers, and S. Das, "Harnessing Voltage Margins for Energy Efficiency in Multicore CPUs," in *Proceedings of the 50th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO '17)*, 2017.
- [76] L. Piga, I. Narayanan, A. Sundarrajan, M. Skach, Q. Deng, M. C. B. Maity, A. Huang, A. Dhanotia, and P. Malani, "Expanding Data-center Capacity with DVFS Boosting: A Safe and Scalable Deployment Experience," in *Proceedings of the 29th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS '24)*, 2024.
- [77] R. P. Potukuchi, J. L. Greathouse, K. Rao, C. Erb, L. Piga, P. G. Voulgaris, and J. Torrellas, "Tangram: Integrated Control of Heterogeneous Computers," in *Proceedings of the 52nd Annual IEEE/ACM International Symposium on Microarchitecture (MICRO '19)*, 2019.
- [78] A. Raghavan, L. Emurian, L. Shao, M. Papaefthymiou, K. P. Pipe, T. F. Wenisch, and M. M. Martin, "Computational Sprinting on a Hardware/Software Testbed," in *Proceedings of the 18th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS '13)*, 2013.

- [79] A. Raghavan, Y. Luo, A. Chandawalla, M. Papaefthymiou, K. P. Pipe, T. F. Wenisch, and M. M. K. Martin, "Computational Sprinting," in *Proceedings of the IEEE 18th International Symposium on High-Performance Comp Architecture (HPCA '12)*, 2012.
- [80] P. Ranganathan, P. Leech, D. Irwin, and J. Chase, "Ensemble-level Power Management for Dense Blade Servers," in *Proceedings of the 33rd Annual International Symposium on Computer Architecture (ISCA '06)*, 2006.
- [81] C. Reiss, A. Tumanov, G. R. Ganger, R. H. Katz, and M. A. Kozuch, "Heterogeneity and Dynamicity of Clouds at Scale: Google Trace Analysis," in *Proceedings of the 3rd Symposium on Cloud Computing (SoCC '12)*, 2012.
- [82] H. Rui, "amd-pstate CPU Performance Scaling Driver," April 2024. [Online]. Available: <https://docs.kernel.org/admin-guide/pm/amd-pstate.html>
- [83] R. Russell, "Virtio: Towards a de-facto standard for virtual i/o devices," *SIGOPS Oper. Syst. Rev.*, vol. 42, no. 5, p. 95–103, jul 2008.
- [84] V. Sakalkar, V. Kontorinis, D. Landhuis, S. Li, D. De Ronde, T. Blooming, A. Ramesh, J. Kennedy, C. Malone, J. Clidas, and P. Ranganathan, "Data Center Power Oversubscription with a Medium Voltage Power Plane and Priority-Aware Capping," in *Proceedings of the 25th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS '20)*, 2020.
- [85] M. Schwarzkopf, A. Konwinski, M. Abd-El-Malek, and J. Wilkes, "Omega: Flexible, Scalable Schedulers for Large Compute Clusters," in *Proceedings of the 8th ACM European Conference on Computer Systems (EuroSys '13)*, 2013.
- [86] M. Shahrad, R. Fonseca, I. Goiri, G. Chaudhry, P. Batum, J. Cooke, E. Laureano, C. Tresness, M. Russinovich, and R. Bianchini, "Serverless in the Wild: Characterizing and Optimizing the Serverless Workload at a Large Cloud Provider," in *Proceedings of the USENIX Annual Technical Conference (USENIX ATC '20)*, 2020.
- [87] J. Stojkovic, N. Iliakopoulou, T. Xu, H. Franke, and J. Torrellas, "EcoFaaS: Rethinking the Design of Serverless Environments for Energy Efficiency," in *Proceedings of the 51st Annual International Symposium on Computer Architecture (ISCA '24)*, 2024.
- [88] J. Stojkovic, T. Xu, H. Franke, and J. Torrellas, "MXFaaS: Resource Sharing in Serverless Environments for Parallelism and Efficiency," in *Proceedings of the 50th Annual International Symposium on Computer Architecture (ISCA '23)*, 2023.
- [89] M. Tirmazi, A. Barker, N. Deng, M. E. Haque, Z. G. Qin, S. Hand, M. Harchol-Balter, and J. Wilkes, "Borg: The next Generation," in *Proceedings of the 15th ACM European Conference on Computer Systems (EuroSys '20)*, 2020.
- [90] A. Verma, L. Pedrosa, M. Korupolu, D. Oppenheimer, E. Tune, and J. Wilkes, "Large-Scale Cluster Management at Google with Borg," in *Proceedings of the 10th ACM European Conference on Computer Systems (EuroSys '15)*, 2015.
- [91] G. Wang, L. Zhang, and W. Xu, "What Can We Learn from Four Years of Data Center Hardware Failures?" in *Proceedings of the 47th Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN '17)*, 2017.
- [92] S. Wang, G. Zhang, J. Wei, Y. Wang, J. Wu, and Q. Luo, "Understanding Silent Data Corruptions in a Large Production CPU Population," in *SOSP*, 2023.
- [93] E. Wu, J. Sune, W. Lai, E. Nowak, J. McKenna, A. Vayshenker, and D. Harmon, "Interplay of voltage and temperature acceleration of oxide breakdown for ultra-thin gate oxides," *Solid-State Electronics*, vol. 46, no. 11, 2002.
- [94] Q. Wu, Q. Deng, L. Ganesh, C.-H. Hsu, Y. Jin, S. Kumar, B. Li, J. Meza, and Y. J. Song, "Dynamo: Facebook's Data Center-Wide Power Management System," in *Proceedings of the 43rd Annual International Symposium on Computer Architecture (ISCA '16)*, 2016.
- [95] Xen Project, "XenStore," April 2024. [Online]. Available: <https://wiki.xenproject.org/wiki/XenStore>
- [96] A. Yassine, H. Nariman, M. McBride, M. Uzer, and K. Olasupo, "Time dependent breakdown of ultrathin gate oxide," *IEEE Transactions on Electron Devices*, vol. 47, no. 7, 2000.
- [97] C. Zhang, A. G. Kumbhare, I. Manousakis, D. Zhang, P. A. Misra, R. Assis, K. Woolcock, N. Mahalingam, B. Warrior, D. Gauthier, L. Kunnath, S. Solomon, O. Morales, M. Fontoura, and R. Bianchini, "Flex: High-Availability Datacenters with Zero Reserved Power," in *Proceedings of the 48th Annual International Symposium on Computer Architecture (ISCA '21)*, 2021.
- [98] J. Zhang, S. Elnikety, S. Zarar, A. Gupta, and S. Garg, "Model-Switching: Dealing with Fluctuating Workloads in Machine-Learning-as-a-Service Systems," in *Proceedings of the 12th USENIX Workshop on Hot Topics in Cloud Computing (HotCloud '20)*, 2020.
- [99] Y. Zhang, W. Hua, Z. Zhou, G. E. Suh, and C. Delimitrou, "Sinan: ML-Based and QoS-Aware Resource Management for Cloud Microservices," in *Proceedings of the 26th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS '21)*, 2021.
- [100] W. Zheng and X. Wang, "Data Center Sprinting: Enabling Computational Sprinting at the Data Center Level," in *Proceedings of the IEEE 35th International Conference on Distributed Computing Systems (ICDCS '15)*, 2015.
- [101] W. Zheng, X. Wang, Y. Ma, C. Li, H. Lin, B. Yao, J. Zhang, and M. Guo, "SprintCon: Controllable and Efficient Computational Sprinting for Data Center Servers," in *Proceedings of the IEEE International Parallel and Distributed Processing Symposium (IPDPS '19)*, 2019.
- [102] L. Zhou, L. N. Bhuyan, and K. K. Ramakrishnan, "Gemini: Learning to Manage CPU Power for Latency-Critical Search Engines," in *Proceedings of the 53rd Annual IEEE/ACM International Symposium on Microarchitecture (MICRO '20)*, 2020.