# Taming Serverless Cold Starts Through OS Co-Design

Ben Holmes, Baltasar Dinis, Lana Honcharuk, Joshua Fried, Adam Belay

MIT CSAIL

## Abstract

Serverless computing promises fine-grained elasticity and operational simplicity, fueling widespread interest from both industry and academia. Yet this promise is undercut by the cold start problem, where invoking a function after a period of inactivity triggers costly initialization before any work can begin. Even with today's high-speed storage, the prevailing view is that achieving sub-millisecond cold starts requires keeping state resident in memory.

This paper challenges that assumption. Our analysis of existing snapshot/restore mechanisms shows that OS-level limitations, not storage speed, are the real barrier to ultra-fast restores from disk. These limitations force current systems to either restore state piecemeal in a costly manner or capture too much state, leading to longer restore times and unpredictable performance. Furthermore, current memory primitives exposed by the OS make it difficult to reliably fetch data into memory and avoid costly runtime page faults.

To overcome these barriers, we present Spice, an execution engine purpose-built for serverless snapshot/restore. Spice integrates directly with the OS to restore kernel state without costly replay and introduces dedicated primitives for restoring memory mappings efficiently and reliably. As a result, Spice delivers near-warm performance on cold restores from disk, reducing latency by up to 14.9× over state-of-the-art process-based systems and 10.6× over VM-based systems. This proves that high performance and memory elasticity no longer need to be a trade-off in serverless computing.

## 1 Introduction

Serverless computing promises effortless elasticity: developers deploy lightweight functions, while the platform dynamically provisions resources in response to demand [16]. Freed from managing infrastructure, users benefit from fine-grained billing and scalability, while providers can, in theory, achieve high utilization by scheduling workloads just-in-time. In today's commercial platforms (e.g., AWS Lambda [1], Azure Functions [2]), these functions are typically packaged in containers or lightweight VMs, which must be initialized before running user code. However, this vision is fundamentally limited by cold starts: the latency incurred when a function is invoked on a machine with no prior cached state [8, 9, 13, 30, 33, 38, 41].

Cold start delays stem from multiple sources, including container setup, language runtime initialization, library loading, and function-specific startup logic such as JIT compilation [20, 34, 40]. These steps often add tens to hundreds of milliseconds — frequently longer than the function's actual
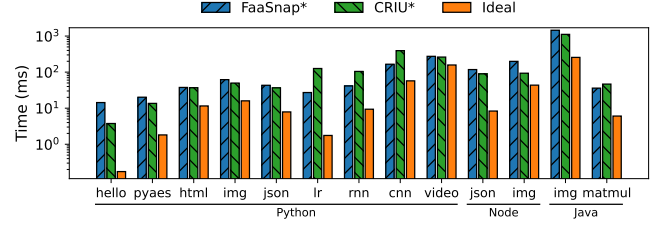


**Figure 1.** The performance gap in snapshot/restore: Even tuned, state-of-the-art systems remain orders of magnitude slower than the ideal restore time. Note the logarithmic y-axis.

execution [10] — making them a major source of user-visible latency. Collectively, these overheads undermine the responsiveness of serverless platforms.

Existing mitigation strategies fall into two broad categories. *Warm-state* approaches, such as keeping containers alive for a short window [4, 15], deliver low-latency startup but consume memory proportional to the number of functions retained. Recent techniques like fork-based replication [7, 10, 17, 41] reduce this cost by enabling replication of a single warm instance across cores or machines. These systems (e.g., leveraging fork locally or using remote memory mechanisms such as RDMA or CXL) allow multiple invocations to share memory copy-on-write, reducing total memory usage while preserving low startup latency. In essence, they amortize one "parent" process across many children, but still require at least one copy of each function to remain alive in memory. This model, however, fundamentally assumes a warm instance is available somewhere on the cluster. For the many functions in the long tail of an invocation distribution, or during a cold start after a period of inactivity, this assumption often does not hold, leaving no fast path for the initial request.

On the other hand, *cold-state* approaches don't rely on warm state present [8, 21, 35, 38]. These approaches often serialize initialized function state to persistent storage and reload it on demand, offering a natural fit for serverless platforms where functions may be invoked rarely but must start quickly. In principle, these snapshot/restore approaches can achieve both elasticity and memory efficiency. In practice, however, restore latencies remain far higher than warm starts, preventing widespread adoption.

To quantify this gap, Figure 1 compares two state-of-the-art systems against a pessimistic *ideal* baseline across several representative serverless workloads (detailed in Table 1). We use an asterisk(*) to denote that we modified both Faasnap

(VM-based) and CRIU (process-based) for optimal performance. For Faasnap, we selected the best-performing of its prefetching strategies, and for CRIU, we eliminated several known overheads. Our ideal baseline is calculated as the sum of the time to read the snapshot's working set from storage and the function's warm execution time, pessimistically assuming no overlap between I/O and computation. Even against these tuned systems and this pessimistic baseline, a significant performance gap remains.

This failure is not due to hardware limits, but to a deeper mismatch between serverless requirements and the interfaces provided by today's OSes. Current OS abstractions are optimized for incremental startup, not for the bulk restoration of an already-initialized process. This mismatch forces existing systems to rely on slow, general-purpose mechanisms. One of the clearest examples is metadata restoration. Lacking dedicated kernel support, process-level tools like CRIU [3] must replay the original setup through a long sequence of expensive system calls. The alternative — snapshotting an entire virtual machine [8, 21, 38] — avoids this replay but captures far more state than necessary, introducing new overheads including scheduling interference from the guest OS. Crucially, both approaches are bottlenecked by a second major challenge: restoring memory contents efficiently.

Memory restoration faces two intertwined challenges. The first is how to proactively populate a process's memory with its working set (the subset of pages needed for execution). One option is to synchronously load the entire predicted working set before execution. While this guarantees the pages are resident in memory, it can result in long delays before execution begins. The alternative is to prefetch asynchronously, overlapping I/O with the function's startup routine. This promises lower latency in theory, but is unreliable with today's OS interfaces and doesn't eliminate stalls due to faults. Second, an application's memory is a composite of shared, file-backed pages and private, modified data. Current OS interfaces cannot restore this complex structure efficiently, forcing a choice between slow, page-by-page updates or sacrificing memory sharing altogether. As fast local storage becomes the norm, it is this two-fold interface mismatch — not raw I/O bandwidth — that has become the dominant factor in cold start latency.

To resolve this interface mismatch and eliminate these OS-level overheads, we present Spice, a snapshot/restore system co-designed with a new set of OS primitives for fast memory and metadata restoration. By providing these missing mechanisms, Spice makes restoring from persistent storage fundamentally practical. Our evaluation shows that Spice reduces the overhead from restoring from a cold disk to under 5ms for a range of functions of varying complexity across several language runtimes.

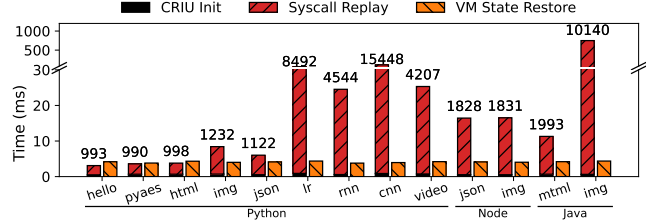This paper makes the following contributions:



**Figure 2.** Initialization overheads for VM- and process-based restore; total syscall counts for replay are shown above the bars. Process restore costs grow with process complexity due to syscall replay. VM restore costs are flat but still non-trivial, reflecting fixed hypervisor operations.

- A study of existing snapshot/restore techniques that identifies persistent overheads due to mismatches between the requirements of low-latency restore and today's OS interfaces.
- A novel OS metadata restoration process that bypasses costly syscall replay by deserializing process state in a single, batched operation from a co-designed snapshot format.
- New kernel mechanisms for memory management that enable fast, reliable restore of process memory that eliminates page fault overhead during startup.
- The implementation and evaluation of Spice, a prototype engine demonstrating that our approach can reduce cold restore latency to under 5ms, substantially outperforming state-of-the-art systems.

Our findings reframe the conventional trade-offs in serverless design, demonstrating that snapshot/restore is a first-class primitive capable of delivering both low latency and high memory efficiency for serverless systems.

## 2 Background and Motivation

This section analyzes the performance of the two predominant snapshot/restore strategies introduced in Section 1: process-level restore, exemplified by CRIU [3], and VM-based snapshots, used by systems like REAP [38], FaaSnap [8], and Sabre [21]. We demonstrate how both approaches are fundamentally constrained by the lack of direct OS support for rapid restoration, leading to the two core bottlenecks foreshadowed earlier: reinstating process metadata and repopulating memory. Our analysis focuses on these core restoration costs; we do not measure the overhead of initializing container primitives like cgroups and network namespaces, as techniques for accelerating these are orthogonal and well-explored [17, 22, 24, 37].

### 2.1 The Challenge of State Reconstruction
The fundamental challenge for fast restoration is that modern operating systems lack a native interface to restore a process's kernel-managed state. This state includes a wide
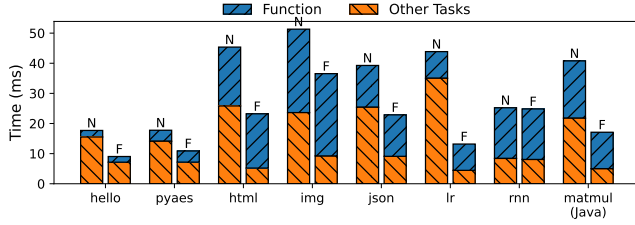
**Figure 3.** Scheduling delay after VM restore with the function running in `SCHED_NORMAL` (N) and `SCHED_FIFO` (F). Even with `SCHED_FIFO`, the function competes with kernel threads and daemons restored from the guest image.

array of resources—threads, file descriptors, memory mappings, signal handlers, and more. Not only is there no single command to reinstate this entire collection at once, but even restoring a single resource to its previous condition can require a complex sequence of operations.

This lack of an efficient restore primitive forces systems into one of two extreme approaches: either (1) replaying the system calls that originally created the process, or (2) snapshotting and restoring the entire virtual machine it runs within. Both strategies push existing abstractions beyond their intended uses and are responsible for significant cold-start overheads.

**Syscall replay.** At one extreme, restore begins from an empty process and attempts to reconstruct every piece of kernel state — threads, file descriptors, memory mappings, and more — by reissuing the original system calls that created them. This strategy is expensive: there are often hundreds or thousands of such calls required, as most kernel objects lack dedicated import or restore interfaces to directly set them to their previous state. For example, restoring a single file descriptor may require a sequence of syscalls: open to create the file, `lseek` to set its offset, dup to assign the correct descriptor number, and `fcntl` to reapply flags. Figure 2 quantifies this overhead, showing that the latency of process-based restore grows with application complexity, from a simple Python function to a more complex JVM application, as the number of syscalls increases from hundreds to thousands. Each resource adds to the latency, and since few can be reconstructed lazily, all must be completed before execution resumes.

The orchestration of this process from userspace adds further overhead. Tools like CRIU stage snapshot data in temporary memory regions at restore time, inject a restorer binary into otherwise unused space, transfer control to that binary, unmap its own code, and finally remap the snapshot data into the correct virtual addresses. This elaborate dance avoids even more costly options such as using a second process to manipulate the target with `ptrace`.

**Full-VM snapshots.** At the other extreme, systems restore entire virtual machines from snapshots. This approach

bypasses syscall replay entirely, as all kernel data structures remain intact. As shown in Figure 2, restore operations are reduced to a small, fixed set of hypervisor actions, such as reinitializing vCPU state, reattaching devices, and restoring other host-managed state. These complete in just a few milliseconds regardless of guest complexity.

A major downside, however, is that reviving the whole VM also revives the entire guest kernel and all its runtime activity — not just the target function. After a prolonged pause, the guest must immediately resume deferred kernel housekeeping and background services. This includes periodic tasks, daemons, and maintenance threads, all of which are scheduled as soon as the VM becomes active. Figure 3 illustrates this problem, showing that a function can be delayed by several milliseconds due to contention with guest kernel tasks like RCU reclamation. We found that even raising the function's scheduling priority with `SCHED_FIFO` offers only partial relief, as critical kernel threads continue to interrupt the function's execution. This scheduling interference significantly impacts end-to-end latency.

*Takeaway:* In the absence of kernel support for reinstating process state directly, systems are pushed into two extremes. Replay-based restore suffers from syscall and orchestration overheads, while VM snapshots avoid replay but revive the entire guest kernel, whose deferred housekeeping and background services cause scheduling interference that heavily impacts function startup latency. Neither approach achieves the millisecond-scale responsiveness required for serverless workloads.

## 2.2 The Challenge of Rapid Memory Restoration

While kernel metadata must be reinstated eagerly, memory contents can be restored lazily. This is attractive in serverless environments where functions often touch only a small portion of their memory. Accordingly, existing systems [8, 21, 38] typically rely on demand paging, bringing pages into memory only when accessed. However, naive demand paging introduces significant latency, as every missing page triggers a blocking fault. To mitigate this, systems employ prefetching to load the predicted working set into memory before it is accessed.

However, as Figure 4 shows, current prefetching strategies are insufficient. REAP [38] employs synchronous prefetching, which loads all pages before execution and minimizes major page faults but results in a long, fixed delay upfront. FaaSnap [8] improves upon REAP's working set estimation and also adds asynchronous prefetching, which attempts to overlap I/O with execution. Unfortunately, this strategy does not meaningfully improve end-to-end latencies, and neither strategy resolves the slowdown incurred by minor faults, caused by missing PTEs or writes to copy-on-write (CoW) pages.

**No Interface for Guaranteed Population.** The fundamental challenge in pre-warming data is the operating system's lack of a reliable, non-blocking interface to populate its page cache. While a synchronous `read()` call can guarantee data is fetched from disk, it does so by stalling the calling thread until the I/O is complete, which results in delayed execution. Consequently, systems must rely on asynchronous, advisory mechanisms through interfaces `madvise()`. However, this is merely a hint, not a command; the kernel retains full discretion to ignore the request, act on it partially, or de-prioritize it based on internal heuristics like memory pressure. This unpredictability means the main application thread remains vulnerable to major page faults on data that was requested but never loaded, ultimately forcing any prefetching strategy into a best-effort approach with no performance guarantees.

**The Lingering Cost of Page Faults.** Even when a page has been successfully fetched into memory, it is not usable until the kernel installs a page table entry (PTE) mapping the virtual address to its physical location. This happens lazily and incurs a *minor fault* on first access, requiring a kernel trap to update the page table. For large heaps or runtime data structures, thousands of such faults can accumulate quickly. Worse, if a fetched page is written to, it will incur an additional fault and copy. File-backed pages are mapped copy-on-write (CoW) to preserve shared page cache state. On the first write, the kernel must allocate a private copy and copy the original contents—an expensive operation. In real-world measurements, JVM-based functions incur tens of thousands of CoW faults during startup, contributing tens of milliseconds of delay.

In virtualized environments, these fault costs become even more severe. Each fault—minor or CoW—forces a VM exit to the hypervisor, incurring hundreds of cycles of latency before fault handling can even begin. Even moderate levels of fault activity can translate into long tail latencies for restored functions running inside VMs like Firecracker.

**Barriers to Efficient Memory Sharing** While prefetching optimizations are important, an orthogonal and equally powerful strategy is to reduce the number of pages that must be fetched in the first place. A natural opportunity arises from reusing memory pages that are already resident in the OS page cache, such as those belonging to shared libraries (e.g., libc) or language runtimes. Figure 5 shows that such shared pages can constitute up to 50% of a function's working set; thus, reusing them could, in principle, halve the volume of data that must be restored from a snapshot.

Existing OS interfaces, however, provide no efficient mechanism for realizing this idea. A process's memory regions, known in Linux as Virtual Memory Areas (VMAs), often contain a mixture of unmodified pages (identical to their backing file) and modified, private pages (unique to the process). The OS tracks this distinction only at the fine granularity of individual page table entries (PTEs), not at the coarser VMA
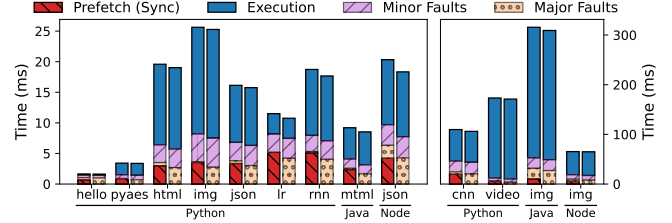


**Figure 4.** Impact of prefetching technique on execution latency. Synchronous prefetching with `read` (left bars) delays execution, while overlapping prefetching with execution using `madvise` (right bars) fails to properly ensure that pages are actually fetched. Neither approach resolves costly minor faults.

level, and it exposes no system call that allows mapping a file while selectively overlaying a small set of private pages.

As a result, restore systems typically fall back on inefficient workarounds. One approach, as implemented in CRIU, is to map the file and then traverse the private pages, issuing a series of costly, one-by-one system calls to update the corresponding PTEs. Another is to fragment what should be a single contiguous memory region into many smaller VMAs in order to isolate modified pages, which inflates kernel metadata and increases memory overhead, and adds to restore latency [40]. Crucially, both strategies require work for every modified range of pages across the application's memory, even for regions that are unlikely to be accessed again, making it impossible to apply overlays lazily or restrict them to the actual working set.

In the context of VM snapshots, the situation is even more restrictive: shared memory reuse across functions is infeasible because the hypervisor treats guest memory and disk as opaque blocks of memory. This opacity prevents the host from identifying and reusing file-backed pages that may already exist in the host's page cache.

*Takeaway:* Prefetching helps avoid reinstating unused memory, but today's kernel interfaces make it difficult to do so efficiently. Taking advantage of cached pages is cumbersome, asynchronous prefetching is unreliable, and minor and CoW faults — even after prefetch — add tens of milliseconds of delay, especially under virtualization. Meaningful improvement will require new kernel mechanisms for precise, non-blocking memory restoration.

## 2.3 Fork is Not a Panacea

Recent work has explored fork-based approaches [7, 10, 17, 41], which achieve near-zero startup latency by cloning a warm process. Within a single machine, these systems leverage the kernel's existing `fork()` mechanism to replicate process state efficiently. This operation is fast because the kernel can duplicate internal data structures in place and relies on copy-on-write to defer actual page copies. In this
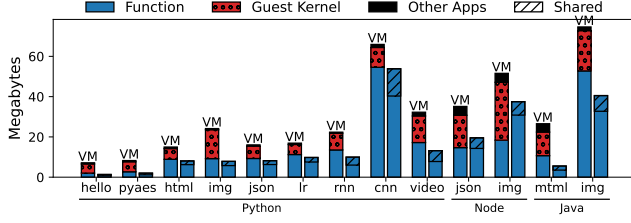
**Figure 5.** Working set composition for VM-based and process-based restore. VMs include substantial kernel-mapped memory and guest-level state, while processes benefit from the potential of sharing 17-51% of working set pages through the host page cache with other functions.



**Figure 6.** Spice Overview

itself a first-class, efficient operation. In the next section, we introduce Spice, which addresses these structural problems directly.

model, metadata restoration is essentially free: file descriptors, mappings, and page tables are inherited wholesale from the parent process rather than reconstructed.

These fork-based approaches are highly effective at rapid, horizontal scaling from a warm parent. Their design, however, makes a distinction between cloning a running instance and restoring that first parent from a cold state. The engineering focus of these systems, particularly for remote fork, is on optimizing high-speed memory replication and cross-node transport. The problem of efficiently instantiating the initial parent from persistent storage, therefore, lies outside their primary optimization path.

This creates a crucial trade-off: forks provide rapid scaling as long as a parent is available, but they do not address the cold-start latency of creating that parent. When all instances have exited, these systems must fall back on slower, conventional methods for instantiation. Our work addresses this specific phase, focusing on an efficient cold restore that can precede the first fork.

### 2.4 Summary of Kernel Limitations

Our analysis reveals that existing restore strategies, whether process- or VM-based, are fundamentally constrained by structural limitations in modern operating systems. The core bottlenecks are twofold: process-based systems are hampered by slow metadata reconstruction, requiring thousands of expensive system calls to rebuild in-kernel state, while both approaches suffer from inefficient memory restoration, relying on unreliable prefetching and a fault-driven process to populate page tables. These issues are compounded by platform-specific overheads, such as scheduling interference within VMs and inefficient memory sharing in container setups. Even fork-based approaches, which excel at cloning warm instances, must ultimately rely on these slower methods when a function is not warm in the cluster.

Ultimately, all current designs expose the same underlying gap: today's OS interfaces are built for incremental startup, not the rapid restoration of a complete process state. Closing this gap requires new kernel mechanisms that make restore
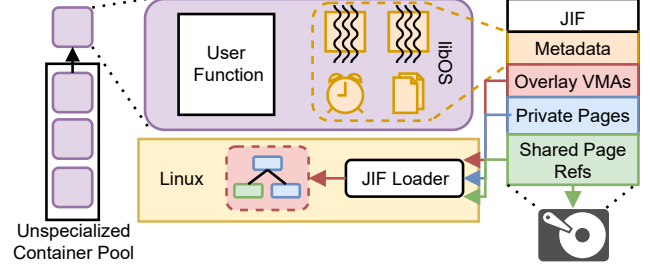
## 3 Approach

Based on the limitations of existing systems, we designed Spice, a new system that leverages novel operating system mechanisms for low-latency function restoration. Our design is driven by the goal of minimizing end-to-end latency while maximizing resource utilization and flexibility. To achieve this balance, our approach is guided by three design goals commonly considered when optimizing cold starts:

1. Minimize the critical-path state by lazily deferring any non-essential initialization.
2. Amortize generic setup costs by performing all function-agnostic preparation ahead of time.
3. Maximize parallelism by overlapping I/O operations with computational work during restoration.

In line with these principles, we choose to operate on process-level snapshots within a container rather than full VMs, as the container is the standard unit of serverless deployment. In addition to reducing the amount of state required to fetch at invocation time (Figure 5), this also provides flexibility: our approach is effective both in bare-metal deployments and within virtual machines, since the function is ultimately executed as a containerized process in either environment. To amortize setup costs, our design utilizes a host-side pool of pre-initialized, unspecialized containers that are ready to be specialized, eliminating container setup costs from the critical path. Our model assumes fast storage is available to load the snapshot into one of these containers quickly.

To realize these ideas, Spice introduces new OS mechanisms to address the core challenges of slow metadata reconstruction and inefficient memory restoration. Effectively addressing these issues requires first-class OS support, as existing interfaces are designed for process creation, not restoration. As a full implementation would demand deep changes to the Linux kernel, we built a prototype in Junction [14], a container system that that implements the Linux kernel interface in userspace (similar to gVisor). Its architecture is ideal for rapid prototyping and offers full control

over the address space layout, which helps avoid the memory mapping conflicts that challenge systems like CRIU. Our ultimate goal is to provide a blueprint for these mechanisms, hoping our results will encourage the Linux kernel community to adopt first-class support for high-performance process restoration. Our prototype design allows us to implement new metadata interfaces entirely in userspace, while our memory management improvements are built as a module in the host kernel. Figure 6 shows our system's design, which we now describe in more detail.

## 3.1 The JIF File Format

To orchestrate the restore process, our approach relies on a new snapshot file format co-designed with our system's architecture. Inspired by the Executable and Linkable Format (ELF) that tells an operating system how to load a program, we introduce the Joint Image Format (JIF). The JIF is a structured binary format that packages all the information needed to restore a process — its serialized kernel metadata, memory pages, and layout information — into a single, self-contained file. This unified format is explicitly designed to be parsed efficiently, enabling Spice to handle metadata and memory restoration in parallel.

## 3.2 Efficient Metadata Restoration

The metadata section of the JIF is generated by a set of intelligent, per-subsystem serialization mechanisms. Instead of treating kernel state as an opaque blob, these interfaces leverage semantic knowledge of each subsystem to create a maximally compact representation. For example, the networking and IPC serializers automatically trim empty pipe and socket buffers, ensuring that only essential in-flight data is saved in the snapshot.

This highly optimized format is key to eliminating the overhead of syscall replay. Spice restores state directly from this compact binary representation, reconstructing the process without replaying individual system calls. Because Junction is a userspace library operating system (libOS), the entire deserialization of this state happens within the Junction runtime, requiring no system calls or kernel transitions. This has two key benefits:

1. No Kernel-Entry Overhead: By handling the entire metadata restore as a single, batched operation within userspace, we avoid the cost of thousands of individual kernel crossings inherent to the syscall replay model.
2. Lazy Resolution and I/O Overlap: The userspace restore process can be fully overlapped with the I/O required to read the snapshot. Where possible, we defer expensive operations. For example, file descriptors are not fully re-opened on restore; instead, they are resolved lazily on their first use, avoiding costly file system traversals during the critical startup path.

## 3.3 High-Performance Memory Management

While Junction manages most of the application's kernel state, the host kernel still manages page tables and memory mappings. Spice introduces a dedicated memory management system in the host kernel with two primary areas of improvement.

### 3.3.1 Optimized Prefetching.
To solve the problems of unreliable, hint-based prefetching and fault-driven page table population, we introduce a new kernel prefetching module. This module leverages the layout of the JIF, which identifies the predicted working set and stores all its constituent pages in a contiguous block. This design is critical for performance, as it enables the module to reliably fetch the entire working set with a single, high-throughput I/O operation, avoiding the overhead of many small reads. After fetching the data, the module pre-populates the corresponding page table entries (PTEs), eliminating the thousands of minor faults that plague existing systems.

The module treats memory differently based on its expected usage to avoid expensive copy-on-write (CoW) faults:

- Pages from the snapshot that were modified during initialization but are not written to during typical execution are mapped as Copy-on-Write, allowing them to be shared safely.
- Pages that are known to be private and written to during restore to are installed directly into writable memory, bypassing the CoW mechanism entirely.

Leveraging the host's page cache provides a powerful mechanism for accelerating frequently invoked functions. By loading pages that are mapped copy-on-write—including both shared libraries and private, read-mostly data—into the cache, the OS can retain this information even after a container is torn down. This enables subsequent invocations to launch rapidly from the warm cache, avoiding the significant resource overhead of explicitly keeping idle container instances alive. This caching benefit does not apply to private, frequently-written pages, whose contents are unique to a single invocation. Spice is designed to handle this distinction, dynamically choosing whether to leverage the page cache for reusable data or restore private memory directly for writable pages based on memory usage patterns.

### 3.3.2 Sharing with Overlay VMAs.
To maximize memory sharing, we introduce a new kernel structure called an Overlay VMA (Figure 7). This allows Spice to efficiently restore a memory region that is mostly shared but contains a sparse set of private, modified pages. The Overlay VMA maps the shared backing file and uses a compact B-tree, stored in the JIF, as an auxiliary data structure. This B-tree serves as a sparse representation of the original page table, indicating which pages are private, modified, or simply zero-filled; this encoding allows us to completely avoid storing or fetching zero pages from the snapshot.
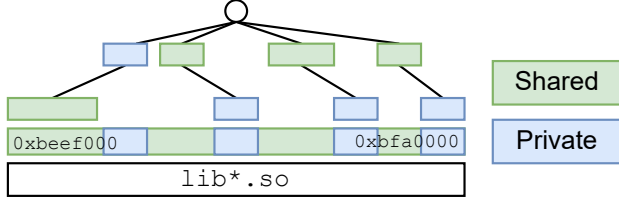
**Figure 7.** An Overlay VMA



**Figure 8.** Spice Detailed Memory Restore Design.

When a page fault occurs within an Overlay VMA, the kernel consults this B-tree to determine how to handle the fault. The tree dictates whether the faulting address should be served from the private overlay in the snapshot, the shared backing file, or by mapping a new zero page. While our optimized prefetching aims to prevent page faults entirely for ideal, profiled workloads, this B-tree mechanism is essential for guaranteeing correctness in case execution diverges. This approach provides a complete solution, avoiding address space fragmentation and costly page-by-page updates without sacrificing correctness.

## 4 Mechanisms for Memory Restoration

Figure 8 details the components and operation of Spice's memory module, which is designed to faithfully recreate the memory subsystem's state as efficiently as possible. Achieving this requires deep integration with memory management to meet its primary goals: prefetching the precise working set to restore memory contents quickly, installing page table entries upfront to avoid page faults, and efficiently preserving memory sharing.

### 4.1 JIF preparation

In an offline phase, functions are pre-warmed with multiple invocations to trigger operations like Just-In-Time (JIT) compilation before a full-memory checkpoint is taken. Memory pages from file-backed mappings are compared with their backing in storage; unmodified pages are discarded from the snapshot, leaving only a reference to the original file. Finally, Overlay VMAs are generated for each original VMA to efficiently represent the overlay of private pages on shared mappings. The Overlay VMA trees are pre-balanced and stored in a compact binary format that requires no deserialization at restore time.

After the initial checkpoint is taken, the system performs several invocations to compute an ordered trace of the function's memory accesses. The trace is added to the JIF file, and the involved pages are relocated and placed in a contiguous range in their order of access. This layout is explicitly optimized to enable high-throughput sequential reading during restoration.
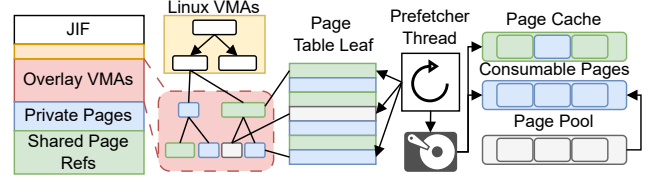
### 4.2 Restore Process

To minimize latency, the memory module is designed to overlap VMA creation with prefetching. These two critical tasks begin the moment a restore is initiated.

Instead of creating and inserting VMAs individually, Spice performs a single batch allocation. This process is accelerated because the pre-computed B-trees for Overlay VMAs are stored in the JIF in a ready-to-use binary format. This allows these complex structures to be directly slotted into the kernel's memory management framework, taking expensive VMA manipulation off the critical path.

Simultaneously, prefetching begins with an initial synchronous I/O batch for the small set of pages the process will access first. This allows execution to resume immediately while the bulk of the memory is fetched by a dedicated prefetcher thread. The prefetcher thread works continuously, interleaving new I/O requests with the installation of PTEs for the previous batch. Eagerly installing PTEs eliminates thousands of minor page faults and ensures pages are ready for the application just before they are needed, preventing execution stalls.

This high-throughput approach is supported by a zero page pool, which is essential for taking page allocation off the critical path. Submitting large batches of I/O requires acquiring many physical pages quickly, which can exhaust fast per-core free lists and fall back to the slower global allocator. The zero page pool provides a ready supply of pre-allocated pages for two purposes: buffering incoming data from the snapshot and satisfying requests for writable zero pages. This prevents the prefetcher from stalling and allows it to keep pace with the application's memory demands.

## 5 Implementation

We implement Spice's metadata snapshot and restore interface within Junction's library OS [14]. Junction provides a containerized environment for running unmodified Linux binaries, offering strong isolation by handling most system calls in userspace within a restrictive seccomp and chroot jail. While we do not use its kernel-bypass features, Junction's userspace management of kernel state is ideal for our purposes. It simplifies metadata serialization and deserialization, allowing us to demonstrate the benefits of a dedicated restore interface. We use the lightweight cereal library to generate our compact metadata representation. To aid in

| Language | Function | Warm Latency ($\mu s$) | Snapshot (Working Set) | | | | | | Description |
|---|---|---|---|---|---|---|---|---|---|
| | | | VMAs | Delta Intervals | Private Pages | Shared Pages | Zero Pages | Working Set (MB) | |
| Java | image | 255,713 | 258 (118) | 708 (253) | 78426 (8378) | 40,951 (1997) | (5867) | (63.4) | Rotate a JPEG |
| | mtml | 5658 | 189 (54) | 549 (99) | 6310 (900) | 40588 (519) | (3) | (5.6) | Matrix Multiplication |
| NodeJS | image | 43,069 | 310 (216) | 566 (330) | 34166 (7894) | 27323 (1708) | (4822) | (56.3) | Rotate a JPEG |
| | json | 8376 | 290 (163) | 469 (200) | 9552 (3655) | 27323 (1347) | (57) | (19.8) | JSON (de)serialization |
| Python | html | 10,977 | 81 (30) | 193 (87) | 3806 (1591) | 3691 (473) | (64) | (8.3) | HTML rendering |
| | cnn | 53,686 | 1912 (176) | 3142 (564) | 72804 (10318) | 367651 (3464) | (413) | (55.4) | CNN inference |
| | hello | 77 | 48 (16) | 154 (56) | 1105 (168) | 2405 (176) | (1) | (1.3) | A no-op function |
| | image | 15,476 | 192 (47) | 328 (97) | 5169 (1471) | 5343 (539) | (10) | (7.9) | Rotate a JPEG |
| | json | 7287 | 154 (40) | 272 (93) | 5002 (1606) | 3999 (492) | (22) | (8.3) | JSON (de)serialization |
| | lr | 1095 | 1292 (98) | 1751 (170) | 25207 (1928) | 39852 (578) | (13) | (9.8) | LR inference |
| | pyaes | 1667 | 49 (16) | 157 (50) | 1240 (302) | 2401 (211) | (1) | (2.0) | AES encryption |
| | rnn | 8634 | 312 (57) | 1437 (164) | 52382 (1537) | 762571 (1013) | (49) | (10.2) | RNN inference |
| | video | 156,751 | 399 (83) | 726 (191) | 7966 (1990) | 56421 (1382) | (1137) | (17.6) | Grayscale conversion |

**Table 1.** Characterization of memory usage for various serverless functions. Delta intervals counts the number of contiguous sets of pages with modified/private application data. Values in parentheses refer to the working set, as opposed to the whole of the snapshot. For zero pages, only those in the working set are reported.

manipulating, updating, and verifying JIF images, we built `jiftool` in 7,400 lines of Rust.

The core memory restoration logic resides in a 4,300-line Linux 6.5.0 kernel module. This module also handles working set estimation. We found that performing this tracing in the kernel is critical for accuracy; the high overhead of userspace tracing tools can stall execution and distort memory access patterns. To generate a stable and accurate profile, our process involves iteratively tracing multiple restores, using the working set from the previous run to pre-populate page table entries (PTEs) for the next, thereby minimizing tracer-induced artifacts.

Our implementation includes a special network-backed file interface to ensure the restored function can begin useful work immediately. Used by per-language shims, this interface allows an invocation request to be queued and made ready before the restore completes. The host OS monitors thread interaction with this interface to identify critical request-handling threads and grants them the highest scheduling priority upon restore, ensuring they execute on their first available time slice.

To minimize snapshot size, we combine application-level cooperation with OS-level optimizations. The network-backed interface provides a channel for the OS to signal an impending checkpoint, allowing the application to proactively perform state trimming. This includes running garbage collection cycles and explicitly dropping pages belonging to freed caches or other transient data structures.

Separately, as a part of any snapshot operation, our OS implementation pursues its own optimizations to minimize unnecessary state. For example, it translates any MADV_FREE calls (lazy memory freeing) into eager MADV_DONTNEED calls to release memory immediately. Additionally, it trims each thread's stack by identifying the current stack pointer and discarding any data in the unused region above its redzone. These complementary techniques ensure the final snapshot is as lean as possible.

## 6  Evaluation

We aim to evaluate Spice by answering the following questions:

1. How does Spice reduce end-to-end cold start latency in the context of existing snapshot/restore systems (§6.1)?
2. How does each technique introduced by Spice contribute to reducing cold start latency (§6.2)?
3. How does Spice perform in response to bursts of invocations (§6.3)?

***Experimental Setup.*** All experiments are conducted on a machine with an Intel(R) Xeon(R) Gold 5420+ with 28 cores @ 2.00GHz and 128 GB of RAM. Our storage device is a Crucial T705 NVMe drive with a max sequential read bandwidth of 13,600MB/s over PCIe Gen 5.0. We run all experiments from this drive and use it to store all checkpoint images and shared libraries used during function invocations. Our test suite includes the memory and CPU intensive functions from FunctionBench [19]. FunctionBench is originally written in Python so we ported two functions each to Node.js and Java to better capture the landscape of serverless functions. Table 1 summarizes the functions used to evaluate Spice.

### 6.1  End-to-End Latency

To understand how Spice performs in the context of existing systems, we compare to existing snapshot/restore systems that restore checkpoints entirely from storage without relying on any warm state. As discussed in Section 2, these systems include VM-based systems that introduce working set prefetching, FaaSnap [8] and REAP [38], as well as CRIU [3] which restores processes entirely from userspace. In light of the observations made in Section 2 we augment the daemon responsible for executing functions in VM-based systems to run in the SCHED_FIFO scheduling class to minimize interference from other tasks running the guest and
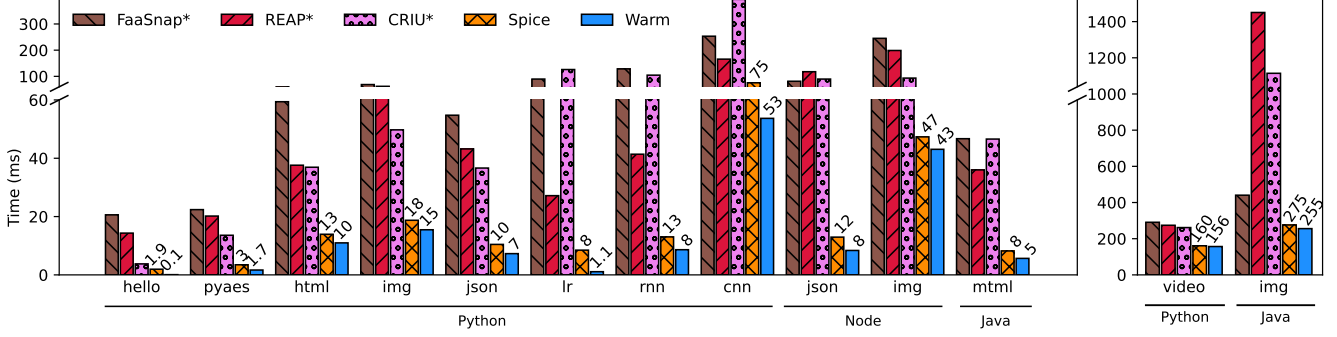
**Figure 9.** Spice achieves end-to-end cold start latencies significantly closer to warm invocations than existing systems.

call them FaaSnap* and REAP*. Native CRIU does not implement working set estimation but instead eagerly installs all memory pages. We modified CRIU to instead use `mmap` and restore memory with demand paging which lowers total latency compared to fetching all of checkpointed memory. We refer to our modified version of CRIU as CRIU*.

Figure 9 shows end-to-end function invocation latency using checkpoints restored entirely from storage with a cold page cache. Spice is able to reduce latency significantly in all cases, by 4-89% compared to REAP* and 7-92% compared to CRIU*. We additionally compare Spice to a warm function invocation. A warm invocation is one of a function that has already been invoked several times but with cold microarchitectural state (CPU caches, TLB). Cold caches illustrate an invocation on an otherwise busy system where cache state has been polluted by other processes running on the machine. With the exception of `hello` which does no computation, Spice is 1.01-7.75× slower than a warm invocation; much of the added cost is due to VMA creation which cannot be parallelized with execution.

Spice is particularly impactful for functions with short execution times, which represent the majority of functions in many serverless deployments [10, 32], but also reduces latency for functions with long execution times. This is a consequence of introducing a dedicated kernel interface – functions with short execution times benefit from optimizations with smaller absolute impact on restore latency, like system call batching and Overlay VMAs, while functions with long execution times benefit from pre-installing PTEs and the page pool. For example, video processing has a large number of anonymous zero pages in its working set and can quickly retrieve a batch of free pages from the page pool for which the prefetcher thread will install PTEs.

### 6.2 Microbenchmarks

**Metadata Restore.** A key component of Spice is the addition of a dedicated interface for restoring OS metadata including threads, VMAs, file descriptor state, signal handlers, and timers. Spice implements a streamlined metadata restore step that leverages the Junction libOS to demonstrate restoration
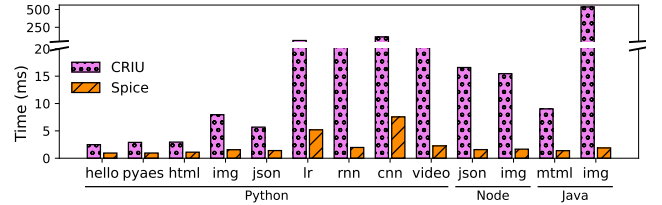


**Figure 10.** OS Metadata restore time comparing CRIU to Spice. Spice lowers latency through a dedicated interface for restoring OS state and compact serialization.

of these components without expensive the system call replay performed in CRIU. In CRIU, metadata restore consists of replaying the system calls that executed during a cold start and expensive de-serialization as user-space data structures are de-serialized and injected into the kernel through system calls and translated into kernel data structures. In Spice, the majority of metadata is held in the userspace data structures of the Junction libOS, while VMAs are recreated in bulk through the kernel module.

Figure 10 compares the metadata restore latency in CRIU to that of Spice, including the time to create VMAs. In all cases, Spice's metadata restore latency is significantly lower than that of CRIU, helping Spice greatly reduce restore latency.

**Memory Restore.** To evaluate Spice's improvements to restoring memory, including actual memory contents, we measure a cold restore through our kernel module with all optimizations disabled and incrementally enable optimizations to illustrate their impact. Figure 11 shows the results with the RNN serving function which performs inference on a small model. Each optimization contributes to an overall large reduction in memory restore time. Overlay VMAs reduce the number of VMAs that need to be created on restore when restoring mappings from shared files that become fragmented when written to. For RNN serving, 1451 VMAs would need to be created to overlay private/modified pages over a shared mapping. The original process has only 314
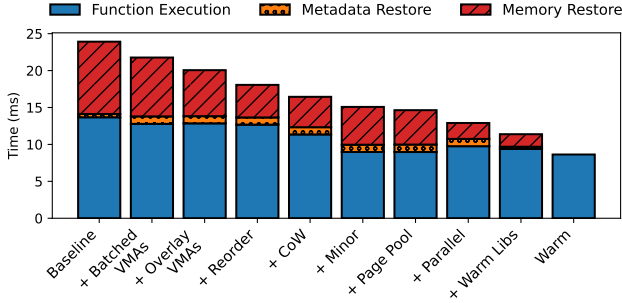
**Figure 11.** Ablation study of Spice's memory restore optimizations on the RNN serving Python function.
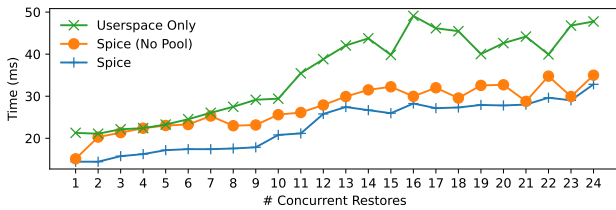


**Figure 12.** Maximum invocation latency as the number of concurrent invocations increases. Spice maintains lower latency compared to an optimized userspace-only baseline due to careful kernel optimizations and lower contention.

VMAs in total. While batched VMA creation reduces the impact of creating a large number of VMAs, Overlay VMAs contribute to reducing memory restore time.

We see additional benefit from reordering the private working set pages in the checkpoint image to be placed in the order they will be accessed by the restored function. Reordering introduces additional VMA fragmentation because adjacent page ranges in the checkpoint file that would be mapped contiguously into a single VMA are split due to accesses that are temporally adjacent but not spatially adjacent. Reordered regions need to be re-shuffled to be adjacent again in the address space, which would necessitate one VMA per region. Reordering in this case creates 3095 total intervals, which we avoid creating through our use of Overlay VMAs.

### 6.3 Concurrency

To see how the techniques introduced by Spice scale with bursts of invocations, we measure the maximum invocation latency in a burst of invocations of a single function from a cold page cache in three configurations: a userspace-only baseline that uses our metadata restore interface but restores memory entirely in userspace, Spice (no pool) uses our memory restore interface with the page pool disabled so private pages are prefetched through the page cache, and Spice with all optimizations enabled. We find that Spice is able to maintain significantly lower latency as the number of concurrent

invocations increases compared to userspace. We also find that the addition of the page pool further reduces latency even though it issues more I/O to read private pages that will be written to during restore. This is likely due to avoiding large page allocations on the critical path, leading to lower latency and more stable performance.

## 7 Related Work

**Remote fork.** Remote fork techniques for function cloning have been explored recently [7, 17, 41]. These systems assume a function is warm in the datacenter and use hardware innovations like RDMA and CXL to quickly fork a new instance on a different machine. We consider these approaches to be orthogonal to ours: in the case that no root function is warm in the datacenter, restoring it from disk is preferable to starting a fresh new instance.

**Sandboxing.** Existing systems use lightweight VMs [4, 8, 15, 21, 38], CRIU [20, 35, 40], or containers [23, 31] for checkpoint/restore. As we discussed in Section 2, these approaches suffer from performance challenges. Some systems use alternative approaches for sandboxing. Faasm [33] relies on Wasm runtime as the isolation mechanism, which offers low startup costs but higher end-to-end execution time compared to a native Linux execution due to the cost of SFI [18]. SEUSS [9] implements unikernels, tailoring their sandbox to executing a specific function. Unlike Spice, SEUSS requires backporting to support additional language runtimes, complicating deployment on existing FaaS platforms.

**Prefetching.** Prefetching the working set of a function to reduce execution time has been explored by prior work [15, 21, 38]. Because there has yet to be a sufficient OS interface for restoring and prefetching function state, these systems suffer from additional overhead with the use of existing solutions. Spice demonstrates that these overheads are not fundamental and other systems could benefit from our design, for example, by snapshotting and restoring a function running in Spice inside a VM if virtualization is desired.

**Avoiding cold starts.** Other systems suggest approaches for pre-warming, optimizing keep-alive policies, and container re-use to avoid cold starts [5, 12, 15, 23, 25–27, 29, 42]. Keeping sandboxes warm hurts resource elasticity, and Spice's fast cold starts represent a step toward eliminating the need for keep-alive policies. Other systems optimize the snapshot timing to ensure that the captured state is the function at its highest performance [9, 44]. Faascale [43] and AFaas [10] both discuss costs associatd with VM EPT faults and pursue strategies to mitigate those.

**Memory Deduplication.** Avoiding memory duplication is key for both performance and resource utilization. Fork-based approaches naturally share memory by leveraging CoW semantics [6, 7, 13, 17, 28, 41]. Other systems like SEUSS [9], Medes [31], and AFaaS [10] propose explicit strategies for de-duplicating snapshot state. SEUSS and AFaaS

create layered snapshot stacks to reuse shared components like language runtimes. However, because these stacks require a strict lineage of changes, their mechanisms are practically confined to sharing a single common base, but not the complex, overlapping mixtures of libraries found across different functions. The Medes system [31] takes a different approach; it reduces the memory penalty of keep-alive containers using explicit deduplication to allow more functions to remain in a semi-warm state. In contrast, Spice aims to reduce the need for keep-alive policies altogether by enabling function state to be stored to and rapidly restored from disk.

**Control path optimization.** While Spice focuses on the data path of function instantiation, significant overheads can also arise from the control path, including request scheduling, resource placement, and network setup. Systems like Dirigent, AFaaS, SigmaOS, and others have explored sophisticated schedulers and resource managers to minimize these orchestration latencies [10, 11, 24, 36, 37, 39]. These efforts are orthogonal and complementary to our work. Achieving the lowest possible end-to-end latency requires a holistic approach, combining an efficient control path with the rapid instance restoration provided by Spice.

## 8 Discussion

Our evaluation demonstrates that by co-designing a snapshot/restore system with new OS primitives, Spice can reduce cold start latency to under 5ms. These results challenge the conventional wisdom that restoring from persistent storage is fundamentally slow. This opens up new possibilities for the design of serverless platforms and raises important questions for future work.

**Resource Management and Keep-Alive Policies.** Spice's sub-5ms restore times in our prototype blur the line between warm and cold starts, potentially altering the economics of serverless resource management. The primary motivation for expensive keep-alive pools—avoiding the high latency of a cold start—is significantly diminished. This enables a new operational model where platforms can practice aggressive reclamation of idle instances to boost utilization, relying on just-in-time instantiation from disk to meet traffic demands without a major latency penalty.

This model also unlocks new strategies for cluster-level optimization. Because Spice leverages the host page cache for sharing, operators can create specialized node pools dedicated to functions with similar software stacks (e.g., a "Python+AI" pool). Co-locating these functions maximizes the page cache hit rate for common runtimes and libraries, which both accelerates restores and increases overall memory density through natural deduplication.

**Integration with Fork-Based Approaches.** The snapshot/restore model of Spice is complementary to the cloning model of fork-based scaling mechanisms [7, 13, 41]. Whereas Spice is optimized for rapidly instantiating the initial "parent" process from persistent storage, fork-based systems excel at cloning that warm parent at microsecond latencies for horizontal scaling.

This relationship suggests a new, hybrid architecture for function instantiation. The principles of Spice's page-cache-aware design could be extended across the network to create a distributed page cache. Such a service would maintain a rack- or cluster-wide pool of frequently-accessed memory pages from common libraries and runtimes. During a restore, Spice could then source required pages from this low-latency remote memory fabric in addition to local storage, further reducing instantiation times and improving resource utilization across the cluster.

**Limitations and Future Work.** A key next step is extending our approach to virtualized environments. While our current techniques can already improve restore times inside a guest VM by reducing in-guest overheads, further potential can be unlocked with direct host-guest cooperation. This could be achieved with a custom hypercall interface that allows the guest to eagerly request EPT population for a dispersed working set. By enabling the guest to pass its predicted memory layout to the hypervisor in a single, batched operation, this approach would eliminate the thousands of costly VM exits typically caused by individual page faults during memory restoration.

Furthermore, while we evaluated with fast local storage, our techniques could be adapted to operate over the network. This would enable restoring functions from remote storage or directly from another node's memory, blurring the line between cold starts and remote fork systems and offering greater placement flexibility in large-scale clusters.

## 9 Conclusion

By demonstrating that cold starts from persistent storage can achieve near-warm latency, this work redefines the fundamental trade-offs in serverless computing. The long-accepted compromise between performance and memory elasticity is not an inherent limitation, but an artifact of operating systems designed for an era before serverless. Our findings suggest that the focus of optimization should shift from user-space heuristics and keep-alive policies to first-class OS support for state restoration.

Spice serves as a blueprint for this new direction. By co-designing the execution engine with the kernel, we unlock a new operational model where functions can be aggressively offloaded to disk to maximize density and efficiency, yet instantiated in milliseconds on demand. This approach opens avenues for future platform architectures built around just-in-time, disk-based instantiation as the default, rather than the exception.

## References

[1] [n. d.]. AWS Lambda. https://aws.amazon.com/pm/lambda/. Accessed: 2025-08-20.

[2] [n. d.]. Azure Functions. https://azure.microsoft.com/en-us/products/functions. Accessed: 2025-08-20.

[3] [n. d.]. CRIU: Checkpoint/Restore In Userspace. https://criu.org/Main_Page.

[4] Alexandru Agache, Marc Brooker, Alexandra Iordache, Anthony Liguori, Rolf Neugebauer, Phil Piwonka, and Diana-Maria Popa. 2020. Firecracker: Lightweight Virtualization for Serverless Applications. In *17th USENIX Symposium on Networked Systems Design and Implementation (NSDI 20)*. USENIX Association, Santa Clara, CA, 419–434. https://www.usenix.org/conference/nsdi20/presentation/agache

[5] Siddharth Agarwal, Maria A. Rodriguez, and Rajkumar Buyya. 2021. A Reinforcement Learning Approach to Reduce Serverless Function Cold Start Frequency. In *2021 IEEE/ACM 21st International Symposium on Cluster, Cloud and Internet Computing (CCGrid)*. 797–803. doi:10.1109/CCGrid51090.2021.00097

[6] Istemi Ekin Akkus, Ruichuan Chen, Ivica Rimac, Manuel Stein, Klaus Satzke, Andre Beck, Paarijaat Aditya, and Volker Hilt. 2018. SAND: Towards High-Performance Serverless Computing. In *2018 USENIX Annual Technical Conference (USENIX ATC 18)*. USENIX Association, Boston, MA, 923–935. https://www.usenix.org/conference/atc18/presentation/akkus

[7] Chloe Alverti, Stratos Psomadakis, Burak Ocalan, Shashwat Jaiswal, Tianyin Xu, and Josep Torrellas. 2025. CXLfork: Fast Remote Fork over CXL Fabrics. In *Proceedings of the 30th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2* (Rotterdam, Netherlands) *(ASPLOS '25)*. Association for Computing Machinery, New York, NY, USA, 210–226. doi:10.1145/3676642.3715988

[8] Lixiang Ao, George Porter, and Geoffrey M. Voelker. 2022. FaaSnap: FaaS made fast using snapshot-based VMs. In *Proceedings of the Seventeenth European Conference on Computer Systems* (Rennes, France) *(EuroSys '22)*. Association for Computing Machinery, New York, NY, USA, 730–746. doi:10.1145/3492321.3524270

[9] James Cadden, Thomas Unger, Yara Awad, Han Dong, Orran Krieger, and Jonathan Appavoo. 2020. SEUSS: skip redundant paths to make serverless fast. In *Proceedings of the Fifteenth European Conference on Computer Systems* (Heraklion, Greece) *(EuroSys '20)*. Association for Computing Machinery, New York, NY, USA, Article 32, 15 pages. doi:10.1145/3342195.3392698

[10] Xiaohu Chai, Tianyu Zhou, Keyang Hu, Jianfeng Tan, Tiwei Bie, Anqi Shen, Dawei Shen, Qi Xing, Shun Song, Tongkai Yang, et al. 2025. Fork in the Road: Reflections and Optimizations for Cold Start Latency in Production Serverless Systems. In *19th USENIX Symposium on Operating Systems Design and Implementation (OSDI 25)*. 199–218.

[11] Lazar Cvetković, François Costa, Mihajlo Djokic, Michal Friedman, and Ana Klimovic. 2024. Dirigent: Lightweight Serverless Orchestration. In *Proceedings of the ACM SIGOPS 30th Symposium on Operating Systems Principles* (Austin, TX, USA) *(SOSP '24)*. Association for Computing Machinery, New York, NY, USA, 369–384. doi:10.1145/3694715.3695966

[12] Nilanjan Daw, Umesh Bellur, and Purushottam Kulkarni. 2020. Xanadu: Mitigating cascading cold starts in serverless function chain deployments. In *Proceedings of the 21st International Middleware Conference* (Delft, Netherlands) *(Middleware '20)*. Association for Computing Machinery, New York, NY, USA, 356–370. doi:10.1145/3423211.3425690

[13] Dong Du, Tianyi Yu, Yubin Xia, Binyu Zang, Guanglu Yan, Chenggang Qin, Qixuan Wu, and Haibo Chen. 2020. Catalyzer: Sub-millisecond Startup for Serverless Computing with Initialization-less Booting. In *Proceedings of the Twenty-Fifth International Conference on Architectural Support for Programming Languages and Operating Systems* (Lausanne, Switzerland) *(ASPLOS '20)*. Association for Computing Machinery, New York, NY, USA, 467–481. doi:10.1145/3373376.3378512

[14] Joshua Fried, Gohar Irfan Chaudhry, Enrique Saurez, Esha Choukse, Inigo Goiri, Sameh Elnikety, Rodrigo Fonseca, and Adam Belay. 2024. Making Kernel Bypass Practical for the Cloud with Junction. In *21st USENIX Symposium on Networked Systems Design and Implementation (NSDI 24)*. USENIX Association, Santa Clara, CA. https://www.usenix.org/conference/nsdi24/presentation/fried

[15] Alexander Fuerst and Prateek Sharma. 2021. FaasCache: keeping serverless computing alive with greedy-dual caching. In *Proceedings of the 26th ACM International Conference on Architectural Support for Programming Languages and Operating Systems* (Virtual, USA) *(ASPLOS '21)*. Association for Computing Machinery, New York, NY, USA, 386–400. doi:10.1145/3445814.3446757

[16] Joseph M. Hellerstein, Jose M. Faleiro, Joseph Gonzalez, Johann Schleier-Smith, Vikram Sreekanti, Alexey Tumanov, and Chenggang Wu. 2019. Serverless Computing: One Step Forward, Two Steps Back. In *9th Biennial Conference on Innovative Data Systems Research, CIDR 2019, Asilomar, CA, USA, January 13-16, 2019, Online Proceedings*. www.cidrdb.org. http://cidrdb.org/cidr2019/papers/p119-hellerstein-cidr19.pdf

[17] Jialiang Huang, MingXing Zhang, Teng Ma, Zheng Liu, Sixing Lin, Kang Chen, Jinlei Jiang, Xia Liao, Yingdi Shan, Ning Zhang, Mengting Lu, Tao Ma, Haifeng Gong, and YongWei Wu. 2024. TrEnv: Transparently Share Serverless Execution Environments Across Different Functions and Nodes. In *Proceedings of the ACM SIGOPS 30th Symposium on Operating Systems Principles* (Austin, TX, USA) *(SOSP '24)*. Association for Computing Machinery, New York, NY, USA, 421–437. doi:10.1145/3694715.3695967

[18] Abhinav Jangda, Bobby Powers, Emery D. Berger, and Arjun Guha. 2019. Not So Fast: Analyzing the Performance of WebAssembly vs. Native Code. In *2019 USENIX Annual Technical Conference (USENIX ATC 19)*. USENIX Association, Renton, WA, 107–120. https://www.usenix.org/conference/atc19/presentation/jangda

[19] Jeongchul Kim and Kyungyong Lee. 2019. FunctionBench: A Suite of Workloads for Serverless Cloud Function Service. In *2019 IEEE 12th International Conference on Cloud Computing (CLOUD)*. 502–504. doi:10.1109/CLOUD.2019.00091

[20] Sumer Kohli, Shreyas Kharbanda, Rodrigo Bruno, Joao Carreira, and Pedro Fonseca. 2024. Pronghorn: Effective Checkpoint Orchestration for Serverless Hot-Starts. In *Proceedings of the Nineteenth European Conference on Computer Systems* (Athens, Greece) *(EuroSys '24)*. Association for Computing Machinery, New York, NY, USA, 298–316. doi:10.1145/3627703.3629556

[21] Nikita Lazarev, Varun Gohil, James Tsai, Andy Anderson, Bhushan Chitlur, Zhiru Zhang, and Christina Delimitrou. 2024. Sabre: Hardware-Accelerated Snapshot Compression for Serverless MicroVMs. In *18th USENIX Symposium on Operating Systems Design and Implementation (OSDI 24)*. USENIX Association, Santa Clara, CA, 1–18. https://www.usenix.org/conference/osdi24/presentation/lazarev

[22] Zijun Li, Jiagan Cheng, Quan Chen, Eryu Guan, Zizheng Bian, Yi Tao, Bin Zha, Qiang Wang, Weidong Han, and Minyi Guo. 2022. RunD: A Lightweight Secure Container Runtime for High-density Deployment and High-concurrency Startup in Serverless Computing. In *2022 USENIX Annual Technical Conference (USENIX ATC 22)*. USENIX Association, Carlsbad, CA, 53–68. https://www.usenix.org/conference/atc22/presentation/li-zijun-rund

[23] Zijun Li, Linsong Guo, Quan Chen, Jiagan Cheng, Chuhao Xu, Deze Zeng, Zhuo Song, Tao Ma, Yong Yang, Chao Li, and Minyi Guo. 2022. Help Rather Than Recycle: Alleviating Cold Startup in Serverless Computing Through Inter-Function Container Sharing. In *2022 USENIX Annual Technical Conference (USENIX ATC 22)*. USENIX Association, Carlsbad, CA, 69–84. https://www.usenix.org/conference/atc22/presentation/li-zijun-help

[24] Zhen Lin, Kao-Feng Hsieh, Yu Sun, Seunghee Shin, and Hui Lu. 2021. FlashCube: Fast Provisioning of Serverless Functions with Streamlined Container Runtimes. In *Proceedings of the 11th Workshop on Programming Languages and Operating Systems* (Virtual Event, Germany) *(PLOS '21)*. Association for Computing Machinery, New York,

NY, USA, 38–45. doi:10.1145/3477113.3487273

[25] David Lion, Adrian Chiu, Hailong Sun, Xin Zhuang, Nikola Grcevski, and Ding Yuan. 2016. Don't Get Caught in the Cold, Warm-up Your JVM: Understand and Eliminate JVM Warm-up Overhead in Data-Parallel Systems. In *12th USENIX Symposium on Operating Systems Design and Implementation (OSDI 16)*. USENIX Association, Savannah, GA, 383–400. https://www.usenix.org/conference/osdi16/technical-sessions/presentation/lion

[26] Wes Lloyd, Minh Vu, Baojia Zhang, Olaf David, and George Leavesley. 2018. Improving Application Migration to Serverless Computing Platforms: Latency Mitigation with Keep-Alive Workloads. In *2018 IEEE/ACM International Conference on Utility and Cloud Computing Companion (UCC Companion)*. 195–200. doi:10.1109/UCC-Companion.2018.00056

[27] Ashraf Mahgoub, Edgardo Barsallo Yi, Karthick Shankar, Sameh Elnikety, Somali Chaterji, and Saurabh Bagchi. 2022. ORION and the Three Rights: Sizing, Bundling, and Prewarming for Serverless DAGs. In *16th USENIX Symposium on Operating Systems Design and Implementation (OSDI 22)*. USENIX Association, Carlsbad, CA, 303–320. https://www.usenix.org/conference/osdi22/presentation/mahgoub

[28] Edward Oakes, Leon Yang, Dennis Zhou, Kevin Houck, Tyler Harter, Andrea Arpaci-Dusseau, and Remzi Arpaci-Dusseau. 2018. SOCK: Rapid Task Provisioning with Serverless-Optimized Containers. In *2018 USENIX Annual Technical Conference (USENIX ATC 18)*. USENIX Association, Boston, MA, 57–70. https://www.usenix.org/conference/atc18/presentation/oakes

[29] Rohan Basu Roy, Tirthak Patel, and Devesh Tiwari. 2022. IceBreaker: warming serverless functions better with heterogeneity. In *Proceedings of the 27th ACM International Conference on Architectural Support for Programming Languages and Operating Systems* (Lausanne, Switzerland) *(ASPLOS '22)*. Association for Computing Machinery, New York, NY, USA, 753–767. doi:10.1145/3503222.3507750

[30] Alireza Sahraei, Soteris Demetriou, Amirali Sobhgol, Haoran Zhang, Abhigna Nagaraja, Neeraj Pathak, Girish Joshi, Carla Souza, Bo Huang, Wyatt Cook, Andrii Golovei, Pradeep Venkat, Andrew Mcfague, Dimitrios Skarlatos, Vipul Patel, Ravinder Thind, Ernesto Gonzalez, Yun Jin, and Chunqiang Tang. 2023. XFaaS: Hyperscale and Low Cost Serverless Functions at Meta. In *Proceedings of the 29th Symposium on Operating Systems Principles* (Koblenz, Germany) *(SOSP '23)*. Association for Computing Machinery, New York, NY, USA, 231–246. doi:10.1145/3600006.3613155

[31] Divyanshu Saxena, Tao Ji, Arjun Singhvi, Junaid Khalid, and Aditya Akella. 2022. Memory deduplication for serverless computing with Medes. In *Proceedings of the Seventeenth European Conference on Computer Systems* (Rennes, France) *(EuroSys '22)*. Association for Computing Machinery, New York, NY, USA, 714–729. doi:10.1145/3492321.3524272

[32] Mohammad Shahrad, Rodrigo Fonseca, Inigo Goiri, Gohar Chaudhry, Paul Batum, Jason Cooke, Eduardo Laureano, Colby Tresness, Mark Russinovich, and Ricardo Bianchini. 2020. Serverless in the Wild: Characterizing and Optimizing the Serverless Workload at a Large Cloud Provider. In *2020 USENIX Annual Technical Conference (USENIX ATC 20)*. USENIX Association, 205–218. https://www.usenix.org/conference/atc20/presentation/shahrad

[33] Simon Shillaker and Peter Pietzuch. 2020. FAASM: lightweight isolation for efficient stateful serverless computing. In *Proceedings of the 2020 USENIX Conference on Usenix Annual Technical Conference (USENIX ATC'20)*. USENIX Association, USA, Article 28, 15 pages.

[34] Wonseok Shin, Wook-Hee Kim, and Changwoo Min. 2022. Fireworks: a fast, efficient, and safe serverless framework using VM-level post-JIT snapshot. In *Proceedings of the Seventeenth European Conference on Computer Systems* (Rennes, France) *(EuroSys '22)*. Association for Computing Machinery, New York, NY, USA, 663–677. doi:10.1145/3492321.3519581

[35] Paulo Silva, Daniel Fireman, and Thiago Emmanuel Pereira. 2020. Prebaking Functions to Warm the Serverless Cold Start. In *Proceedings of the 21st International Middleware Conference* (Delft, Netherlands) *(Middleware '20)*. Association for Computing Machinery, New York, NY, USA, 1–13. doi:10.1145/3423211.3425682

[36] Arjun Singhvi, Arjun Balasubramanian, Kevin Houck, Mohammed Danish Shaikh, Shivaram Venkataraman, and Aditya Akella. 2021. Atoll: A Scalable Low-Latency Serverless Platform. In *Proceedings of the ACM Symposium on Cloud Computing* (Seattle, WA, USA) *(SoCC '21)*. Association for Computing Machinery, New York, NY, USA, 138–152. doi:10.1145/3472883.3486981

[37] Ariel Szekely, Adam Belay, Robert Morris, and M. Frans Kaashoek. 2024. Unifying serverless and microservice workloads with SigmaOS. In *Proceedings of the ACM SIGOPS 30th Symposium on Operating Systems Principles* (Austin, TX, USA) *(SOSP '24)*. Association for Computing Machinery, New York, NY, USA, 385–402. doi:10.1145/3694715.3695947

[38] Dmitrii Ustiugov, Plamen Petrov, Marios Kogias, Edouard Bugnion, and Boris Grot. 2021. Benchmarking, analysis, and optimization of serverless function snapshots. In *Proceedings of the 26th ACM International Conference on Architectural Support for Programming Languages and Operating Systems* (Virtual, USA) *(ASPLOS '21)*. Association for Computing Machinery, New York, NY, USA, 559–572. doi:10.1145/3445814.3446714

[39] Ao Wang, Shuai Chang, Huangshi Tian, Hongqi Wang, Haoran Yang, Huiba Li, Rui Du, and Yue Cheng. 2021. FaaSNet: Scalable and Fast Provisioning of Custom Serverless Container Runtimes at Alibaba Cloud Function Compute. In *2021 USENIX Annual Technical Conference (USENIX ATC 21)*. USENIX Association, 443–457. https://www.usenix.org/conference/atc21/presentation/wang-ao

[40] Kai-Ting Amy Wang, Rayson Ho, and Peng Wu. 2019. Replayable Execution Optimized for Page Sharing for a Managed Runtime Environment. In *Proceedings of the Fourteenth EuroSys Conference 2019* (Dresden, Germany) *(EuroSys '19)*. Association for Computing Machinery, New York, NY, USA, Article 39, 16 pages. doi:10.1145/3302424.3303978

[41] Xingda Wei, Fangming Lu, Tianxia Wang, Jinyu Gu, Yuhan Yang, Rong Chen, and Haibo Chen. 2023. No Provisioned Concurrency: Fast RDMA-codesigned Remote Fork for Serverless Computing. In *17th USENIX Symposium on Operating Systems Design and Implementation (OSDI 23)*. USENIX Association, Boston, MA, 497–517. https://www.usenix.org/conference/osdi23/presentation/wei-rdma

[42] Hanfei Yu, Rohan Basu Roy, Christian Fontenot, Devesh Tiwari, Jian Li, Hong Zhang, Hao Wang, and Seung-Jong Park. 2024. Rainbow-Cake: Mitigating Cold-starts in Serverless with Layer-wise Container Caching and Sharing. In *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 1* (La Jolla, CA, USA) *(ASPLOS '24)*. Association for Computing Machinery, New York, NY, USA, 335–350. doi:10.1145/3617232.3624871

[43] Xinmin Zhang, Qiang He, Hao Fan, and Song Wu. 2024. Faascale: Scaling MicroVM Vertically for Serverless Computing with Memory Elasticity. In *Proceedings of the 2024 ACM Symposium on Cloud Computing* (Redmond, WA, USA) *(SoCC '24)*. Association for Computing Machinery, New York, NY, USA, 196–212. doi:10.1145/3698038.3698512

[44] Yifei Zhang, Tianxiao Gu, Xiaolin Zheng, Lei Yu, Wei Kuai, and Sanhong Li. 2021. Towards a Serverless Java Runtime. In *2021 36th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. 1156–1160. doi:10.1109/ASE51524.2021.9678709