

# Asynchronous I/O — With Great Power Comes Great Responsibility

Constantin Pestka  
*German Aerospace Center*

Marcus Paradies  
*LMU Munich*

Matthias Pohl  
*German Aerospace Center*

## Abstract

The performance of storage hardware has improved vastly recently, leaving the traditional I/O stack incapable of exploiting these gains due to increasingly large relative overheads. Newer asynchronous I/O APIs, such as `io_uring`, have significantly improved performance by reducing such overheads, but exhibit limited adoption in practice. In this paper, we discuss the complexities that the usage of these contemporary I/O APIs introduces to applications, which we believe are mostly responsible for their low adoption rate. Finally, we share implications and trade offs made by architectures that may be used to integrate asynchronous I/O into DB applications.

## 1 Introduction

SSDs have witnessed substantial performance advancements in recent years. Even commodity SSDs can achieve low double-digit  $\mu$ s latency, double-digit GiB/s throughput, and can process millions of IOPS [13, 19, 21]. These performance advances pose tremendous challenges to contemporary I/O software stacks. Overheads, such as frequent context switches in blocking I/O, have made it increasingly difficult to fully utilize modern hardware efficiently. Consequently, there has been a variety of novel I/O APIs that attempt to address these performance issues, i.e., NVMe [2] at the protocol level, `io_uring` [4] and I/O Rings [1] at the OS level, SPDK [24] for OS bypassing access, and `xNVMe` [17] as a unifying abstraction layer on top. One unifying design aspect of these APIs is the transition from a blocking to an asynchronous, completion-based interface. Depending on the used I/O API and configuration, this communication mechanism largely replaces the need for syscalls with mostly lock-less inter-thread communication between multiple kernel and user-space threads. Thus, these new async I/O APIs have the potential to be significantly more efficient than their blocking counterparts [4, 6, 11, 16, 20]. However, shifting to these APIs involves moving from a blocking, preemptive multitasking design to an explicitly async, cooperative multitasking approach. This change transfers the

responsibility for managing the time between I/O request submission and receiving the results to the application. This may not pose an issue for applications that are already multi-threaded, e.g., by using an async job system or an explicit event loop. However, it may require significant refactoring for simpler multi-threaded or single-threaded applications.

While it is possible to issue low MIOPS per thread (cf. Figure 1, [4]), the sheer extent to which the throughput capability of I/O devices has grown in recent years, vastly outpaces the capability of a single thread on essentially all contemporary platforms. Current server generations are frequently equipped with up to 128 PCIe 5.0 lanes that can accommodate dozens of SSDs, each capable of serving multiple MIOPS [3]. Achieving a saturation of 10-20% on such a system fundamentally requires multiple dedicated threads worth of CPU time. Applications that wish to fully utilize modern I/O hardware thus have to implement some degree of parallelism on top of these APIs. Unfortunately, implementing parallelism is generally non-trivial (§ 3, 4), especially while retaining acceptable levels of performance. Despite the significant performance potential the adoption of modern async I/O APIs in I/O heavy applications is low. Many state-of-the-art I/O libraries, such as `libuv`, `seastar` and `tokio` [5, 8, 14] have only limited experimental support and widely used libraries, such as `libc` and `libc++`, do not support novel I/O APIs such as `io_uring` at all. The few libraries that do support these, integrate them in a relatively limited, experimental fashion [7, 8]. Similarly, while async I/O APIs have been explored in academic DBMSs [11, 23] and commercial DBMSs [9], the majority of large, commercial DBMSs still rely on blocking I/O APIs. We believe this is due to the many, often complex responsibilities that these APIs bring to applications (e.g., the requirement to implement efficient parallelism and user-space task scheduling).

In this paper we will first in § 2 briefly present an overview of the benefits of async I/O and will then in § 3 in detail go over the not so frequently discussed challenges users of async I/O have to cope with. In § 4, we conclude with a discussion on a set of architectural patterns that can be applied

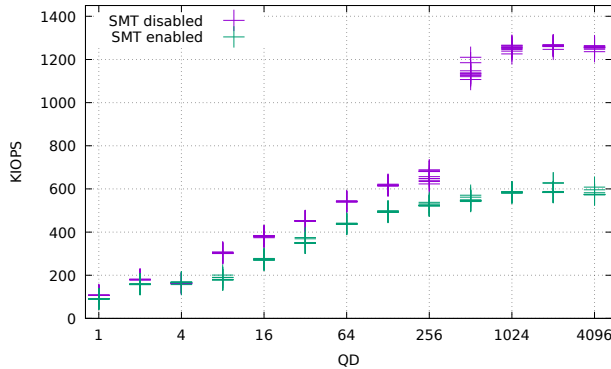


Figure 1: IOPS queue depth scaling. Uses `io_uring`, direct I/O, SQ and I/O poll, 1M 4K sequential reads per run, 10 runs each, preceded by preconditioning run. Run on a Ryzen 3900X, Samsung 990 Pro 1TB and Linux 6.8.0-39.

for the purpose of using async I/O in applications and how they relate to the previously discussed challenges. In this discussion, we will include classifications of existing systems that already utilize async I/O and discuss trade-offs made in their implementations.

## 2 Background

The design of blocking I/O APIs dates back to a time when storage was slow and mostly based on magnetic media. Multiple decades of CPU and storage hardware evolution have continuously eroded the initial assumptions made during the design of traditional blocking I/O APIs. On the CPU side, single-core performance improvements have stagnated, as its main historical driver of frequency scaling has slowed notably. Current CPU performance has mostly improved due to IPC gains and is largely limited by the memory subsystem, which itself heavily depends on its caching hierarchy. This is particularly critical for context switches, as caches have to be flushed to avoid side channels and thus additionally cause high indirect costs [10, 18, 22, 25]. On the storage side, modern SSDs, while experiencing similar frequency scaling issues and thus exhibiting lower gains in latency, provide massive I/O throughput by exploiting a high degree of internal parallelism. In the following we provide an overview of the challenges imposed by blocking I/O on modern hardware and how async APIs attempt to address these.

### 2.1 The Inefficiencies of Blocking I/O

Blocking I/O APIs delegate the responsibility to schedule work between I/O request submission & completion to the OS rather than the application. To access the OS scheduler used for this purpose and to access the kernel I/O stack, user space threads have to invoke a syscall to submit any I/O operation. The resulting, at the very least two, context switches are a the fundamental performance issue of blocking I/O APIs.

Another issue is that the process, which issues a blocking I/O request, can only make further forward progress if it poses an additional, currently not running, not already blocked thread, to have a *chance* to be scheduled to run upon submitting the blocking I/O request. Further, while depending on the configuration, the process is often more likely to do so the more runnable threads it poses. This can lead to over-subscription, which can be disadvantageous for overall system performance and efficiency. While heavily depending on configuration knobs, e.g., thread priorities, scheduling classes & cgroup configuration, the mostly stochastic nature of the OS scheduling can also be undesirable as it might lead to unpredictable behavior. Over-subscription can be problematic in the context of modern SSDs, as exploiting their vast bandwidth is predicated on a sufficiently high number of concurrently in-flight requests (cf. Figure 1). To achieve this an application can either utilize additional threads or use vectored I/O. Vectored I/O allows the submission of multiple I/O requests via a single syscall. While this reduces the required syscalls per I/O operation, it does not allow for different I/O request types to be mixed and requires all requests to be known at the time of submission. The latter can be problematic in applications where the need for I/O requests does not occur in large batches or is unpredictable. Further, while in the multi-threaded approach, the application can always immediately utilize the result of every I/O request once it completes. In the vectored approach, the thread generally remains blocked until all results arrive.

### 2.2 The Boons of Async I/O

Most I/O operations are fundamentally asynchronous. Blocking I/O APIs try to hide this fundamentally async nature and the accompanying complexity from the user. Instead, async I/O APIs fully expose this to the user. This splits the, for the user, essentially atomic I/O operation of blocking APIs into two distinct events, the submission and completion, and moves the scheduling of tasks for the interim into user space, allowing most syscalls be eliminated. While some old I/O APIs (e.g., `aio`) have not taken this opportunity and consequently have been criticized [4], newer I/O APIs (e.g., `io_uring`, I/O Rings, or SPDK) have done so. They typically achieve this via two lockless ring buffers, the *Submission Queue* (SQ) through which I/O requests are submitted, and the *Completion Queue* through which notifications of completion are received. These two communication channels can vastly reduce or completely alleviate the need for syscalls and replace them with lower overhead atomic memory operations. The exact details of when syscalls are still required highly depend on configuration parameters, such as using *submission queue* and *I/O device polling* (SQ and IO poll), but have significant impact on performance [4, 6, 20]. Another advantage of async APIs is that applications can utilize the higher degree of available context to make more optimal scheduling decisions.

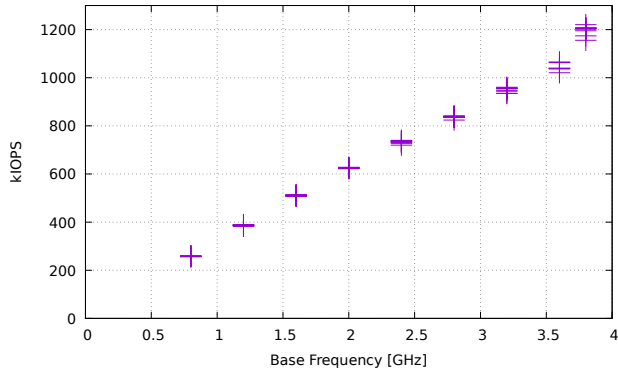


Figure 2: Frequency set in the BIOS (Parameters cf. Figure 1)

Instead of relying on the coarse per thread prioritization schemes of the OS scheduler, applications may optimize their scheduling for generally higher throughput, latency or QoS or choose to prioritize specific groups of tasks, such as persisting critical logs. For SSDs the most critical benefit is the improved efficiency of achieving a high number of concurrent IOPS. Here, async APIs combine the advantages of the threaded, blocking approach and vectored I/O, while not suffering from their limitations. Results of individual requests can always be utilized upon arrival and are not limited by batching. Further amortizing costs via batching is possible, but failing to do so is significantly less costly as the overall submission overhead is vastly reduced.

### 3 The Banes of Async I/O

While, async I/O is conceptually rather simple and provides many opportunities to improve the efficiency of I/O-intensive applications, actually doing so is anything but trivial for all but the simplest applications. Utilizing these APIs, especially when scaling parallelism to fully utilize modern hardware, introduces complex challenges that we believe have not been sufficiently addressed so far.

#### 3.1 I/O has become CPU bound

While historically I/O performance has been limited by I/O device speeds, recently it has become increasingly CPU-limited, and lead to the development of novel async I/O APIs. And while these APIs have massively improved CPU efficiency and hence also maximally achievable performance on one thread [4, 6, 11, 16, 20], even with these APIs, I/O performance on one thread is still closely linked to CPU performance (cf. Figure 2) The CPU time of a single thread is not sufficient to saturate modest modern hardware setups, such as a single PCIe 5.0 SSD, even in idealized benchmark scenarios. The first implication of the degree to which I/O performance is CPU-limited on modern hardware is that a significant degree of parallelism via multi-threading is strictly required to saturate contemporary hardware setups. The second implication

is that the budget of CPU time both for handling I/O related tasks and any other tasks, is limited. Even short delays, e.g., executing other application logic, processing the result of finished IOPs or a few LLC misses, can reduce performance by orders of magnitude (cf. Figure 3). Techniques to retain high performance on modern CPUs & memory subsystems, include utilizing small, non-fragmented, cache-efficient objects & data structures and custom memory allocators. Although they are in heavy use in some communities, such as the Linux kernel, they are not yet commonly used in I/O heavy applications. Additionally, while SSD performance characteristics themselves are non-trivial due to internal factors (e.g., page size, garbage collection, SLC caching), on top of these, the factors that impact general CPU performance now also impact I/O performance. This includes software induced factors, e.g. cache misses, but also setup related factors, such as dynamic frequency scaling, which is strongly dependent on thermal headroom and the power management settings. Among these, we found that SMT has a pronounced impact on performance (cf. Figure 1). On modern hardware these factors have to be considered both for I/O heavy applications and in I/O related benchmarks, not only due to their impact on performance, but also due to the associated understandability and reproducibility issues.

#### 3.2 Moving Scheduling to User Space

While we have previously touched on the advantages of moving scheduling responsibilities to user space, this also has notable downsides. Scheduling itself is a famously non-trivial topic. This, especially in combination with the performance requirements and necessity of multi-threading, introduces an additional burden on the application implementations. This is only complicated further, as the application has to be able to efficiently cope with the async nature of the APIs, e.g., via async job systems that run on top of a thread pool (§ 4.2). Furthermore, the discontinuous nature of the logical control flow of async tasks can complicate debugging, as often a full callstack does not exist.

#### 3.3 Per API Instance Overheads

One of the significant differences between OS-level and OS-bypassing APIs is that instances of the former can come with a steep upkeep cost, most notably due to SQ polling kernel threads. Similar to I/O polling, which removes the need for expensive IRQs to handle arriving results at the cost of additional CPU resources, SQ polling removes the need for syscalls during submission at the cost of a dedicated thread per API instance. While optional, both polling modes significantly increase performance [4, 6, 20] at the cost of the additional CPU time. As these threads are implicitly controlled via a *ms* granularity timeout, they remain active even if only few requests are submitted to them in small enough time in-

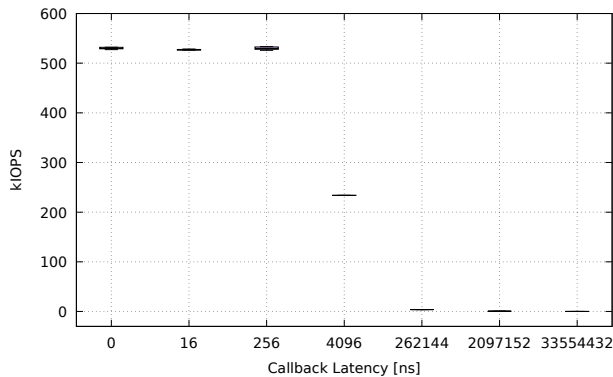


Figure 3: IOPS scaling with post I/O operation callback latency (Parameters cf. Figure 1, but using random reads)

tervals. This can be problematic in architectures, where it is difficult or impossible to match the number of API instances to the hardware capabilities. Failing to adapt the active number of API instances to a varying I/O load can, depending on the workload, be similarly detrimental. The significant upkeep cost of the polling threads generally imposes the challenge for architectures to balance the overheads introduced by sharing API instances with the cost of low API instance utilization.

### 3.4 Task Division and Task Distribution

Two key architectural decisions are how logical tasks are divided into subtasks (§ 4.1) and how the execution of these subtasks is organized among threads (§ 4.2). There is a trade-off between schemes that divide at a smaller granularity increasing dispatch overheads, such as function calls, persisting metadata of requests while they are in flight, and potentially inter-thread communication and schemes dividing at a larger granularity, which makes utilizing dedicated I/O threads (§ 3.5) impossible. Similarly, the scheme chosen to distribute requests among threads has to balance the amount of inter-thread communication and the utilization of API instances.

### 3.5 To I/O Thread or not to I/O Thread

It is a consequential question whether or not all I/O-related tasks, such as the submission & polling for completions, managing the metadata of in-flight requests, and handling partial completions, should be isolated into separated subtasks that are then executed on dedicated *I/O threads*. Not doing so has the advantage that it avoids the need for additional inter-thread communication. However, executing any non-I/O related subtask on the same thread that submits I/O (*inline execution*) will reduce the API instance utilization to a degree proportional to the execution time of the non-I/O tasks executed. Given the extent to which throughput and hence API instance utilization drops based on the execution time of tasks executed inline (cf. Figure 3), CPU efficiency will drop significantly if inline execution is used alongside I/O or SQ polling, even for

short non-I/O tasks of a few  $\mu$ s.

## 3.6 Increased complexity over traditional I/O

Async APIs (e.g., `io_uring`) are more complex than blocking I/O, as they include more features and introduce additional responsibilities to the user. In the following we discuss some of them. As the SQ and CQ of `io_uring`'s API instances are *statically* sized ring buffers, the application has to avoid overflowing the SQ. Enforcing ordering and atomicity among I/O operations is more involved, especially if cancellations are involved, which also makes the semantics of `fsync` more complex than usual. `io_uring` exposes features, such as multi-shot accepts, registered file descriptors & buffers, which can improve efficiency, but also increase complexity. Finally, to correctly use an API instance with more than one thread, familiarity with modern memory models is required to achieve correctness. Furthermore, most programmers use high-level libraries, such as `libc` to perform I/O. However, the vast majority of I/O libraries do not support the new async I/O APIs and the few that do (e.g., [7, 8]) do so in a limited fashion. While an exhaustive discussion of these libraries is out of the scope of this paper, limitations such as utilizing a single-threaded event loop and, thus, limited throughput limit their scope of usability on modern hardware. The complexities of a correct and efficient implementation that utilizes e.g. `io_uring` is thus highly relevant, as most applications that wish to utilize these more efficient I/O APIs will have to implement the integration on their own.

## 4 Async I/O Architectures

### 4.1 Task Partitioning Schemes

There are three common ways to partition a task including async I/O. To compare these, it is helpful to divide a task into subtasks separated by I/O operations. We then define a *tasklet* as a sequence of instructions including one or more sub-tasks that one scheduled run from start to finish without interruption on the same thread.

#### 4.1.1 Full Partitioning

In *Full Partitioning*, every subtask is divided into its own tasklet. If there is a subsequent IOP, then the submission of the according request is appended to the tasklet of the previous subtask. The polling for completion of an IOP is encapsulated into a separate tasklet, which spawns itself again on failure and the tasklet of its successor subtask on success. This approach is highly flexible, as it allows implementations to choose any distribution among threads and prioritization scheme.

### 4.1.2 Callback Partitioning

Another method similar to full partitioning is to use a callback function that is called after the polling of an I/O request returns a successful status. This method combines the tasklet for polling I/O requests with the subsequent subtask. It should be noted that this implies that said sub-task is also executed by the same thread that performed the polling of the I/O request, hence preventing the usage of I/O threads, which is not the case for full partitioning.

### 4.1.3 Coroutines

Coroutines are functions that can be suspended and resumed at predefined points, here most notably after IO submission and unsuccessful polling for completion and thus can be used to encapsulate an entire task. They are typically implemented as stackless coroutines, which consist of a heap-allocated stackframe containing the function arguments, current locals, and a state enum. When a coroutine is called the state enum is used to determine the point at which execution is to be resumed. Some implementations enforce a specific execution architecture, such as a single-threaded eventloop, thus also always introducing the limitations of said architecture, such as limited parallelism. The properties of coroutines thus vary wildly between implementations of different languages such as C++, Rust, JavaScript, Zig, Go or manual implementations [15] and are hence difficult to discuss in this context. One notable performance implication of stackless coroutines is that resuming coroutines, especially nested ones, can be rather expensive, as coroutines frames are heap allocated and often large and thus hostile to caches, as they must fit the entire current state of the full task. As polling for I/O completions must remain inexpensive, using coroutines for async I/O can be problematic. A notable upside of coroutines is that they encapsulate more context of the current task in the callstack as they contain the full task, which can be beneficial for debugging. Coroutines thus have been a popular choice for async I/O [23] or other async operations like hiding pre-fetching operations [12, 15].

## 4.2 Execution Architectures

In this chapter, we discuss the distribution of tasklets among threads for execution. Our discussion of architectures will focus on their ability to scale efficiently to higher degrees of parallelism. Due to this we will exclude fully single threaded applications and single threaded event loops, that can be encountered in e.g. the Zig standard library or the JS runtime, from our discussion, as their limited throughput is insufficient for modern hardware. Similarly, we will limit the discussion to *thread per core* architectures and thread per core compatible architectures, as the alternative approach of over-subscription and preemptive multi-tasking are problematic in this con-

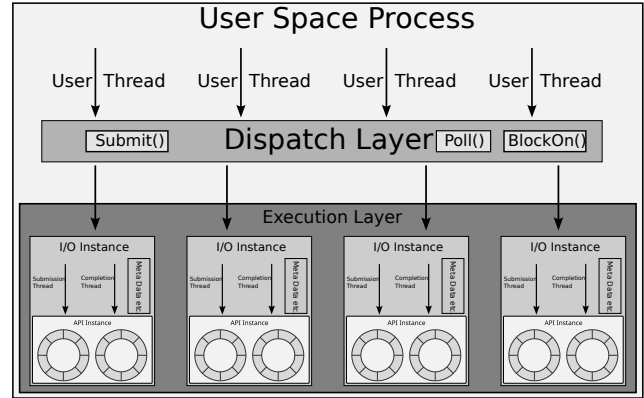


Figure 4: Static I/O thread pool

text due to their often excessive amount of context switches, whose problematic nature has been discussed previously.

### 4.2.1 Direct Access or M-to-N sharing

The significant upkeep cost of API instances is a strong incentive to utilize them efficiently. A simple strategy to achieve this is to share them between multiple threads, e.g. by allowing all of the user space threads ( $N$ ) of an application to access any of the API instances ( $M$ ) directly. This has the advantage that the resources allocated for I/O and non-I/O tasks can be adjusted by changing  $M$  and  $N$ , respectively. While this mostly decouples the amount of resources allocated for I/O from the ones allocated for other purposes, as the CPU time to perform e.g. submissions is still provided by the  $N$  user space threads, the chosen configuration is static and cannot be adapted later. In scenarios with a low I/O load, the submission of requests can thus not be dynamically restricted to a smaller set of API instances resulting in low CPU efficiency (§ 3.3). Additionally, this architecture does not provide the capability to buffer requests if a thread attempts to submit to an API instance with a full SQ and thus must provide this itself. Additionally, it is the responsibility of each thread to choose an API instance to submit to, which could cause congestion issues. The most relevant property of this architecture for our purposes is, however, its issues with scalability. Not only are more expensive synchronization primitives required compared to architectures that restrict access to a specific API instance to only one dedicated thread, but more importantly, the potential for contention is significantly higher by design than for the other presented architectures. This is especially relevant for the polling operations on the CQs. Likely due to said performance issues we are not aware of any project utilizing this architecture.

### 4.2.2 Shared Nothing

In this architecture, each user-space thread submits any of its I/O requests to its private API instance(s). This architecture proposes to achieve parallelism via multiple fully independent

instances of the single-threaded architecture. This approach has a few significant upsides, in particular its simplicity and requiring *no* inter-thread communication.

However, this architecture is not applicable in its pure form in most real-world applications. While it works remarkably well in benchmarks, where requests are generated right before submission, and results are discarded immediately upon completion, this effectively ignores that in most real applications, there is logic to be executed that causes an I/O operation and processes its result. In most applications chains of logic operations, often triggered by non-compile-time known events such as read input from a file or received network packet, eventually require one or more I/O OPs to be issued, of which the output is then processed further. This architecture assumes that the workload *can* be efficiently partitioned into tasks that are then executed fully independently from each other. In other words, this assumes the task is at least approximately embarrassingly parallel and can be efficiently statically scheduled. The effectively non-deterministic latency of I/O, especially for network I/O, conflicts with this assumption. Most I/O intensive applications, like DBMSs or web servers are inherently dynamic, with numerous compile-time unknowns and often complex dependencies. The resulting complexity and interconnectedness of their control flow graphs are fundamentally at odds with being easily and efficiently partitionable. That not all or even most tasks are embarrassingly parallel is a well-known fact in many domains, such as GPGPUs.

Another property of this architecture is that, due to not communicating among threads, I/O throughput available to a task is limited to the maximum throughput of one thread. Further implications are due to utilizing inline execution (§ 3.5), especially regarding CPU efficiency. This is especially critical in this architecture, as there is a hard coupling between available resources for non-I/O and I/O tasks. This is due to the fact that every user space thread must have at least one dedicated API instance and thus, if SQ polling is used, comes with a dedicated polling thread. This severely restricts the maximum amount available resources for non-I/O tasks, while not resorting to over-subscription, as e.g. 32 user space threads, that are also still required to perform the I/O related work, submission, polling etc., would necessitate 32 kernel threads even when using the minimum of one API instance per thread.

### 4.2.3 Static I/O Thread pool

This architecture (cf. Figure 4) attempts to address both, the inter-thread synchronization-based performance issues of the direct access architecture and the restrictions with respect to the applicability of the shared-nothing architecture. To achieve this, it adds a dispatch layer to which all user-space threads submit requests and receive an object from that can be polled for the status of the submitted request. The distribution layer distributes the submitted requests to one of potentially multiple *I/O instances*. An *I/O instance* consists of either two

threads, one responsible for submission one for completion, or one thread responsible for both, as well as one or more I/O API instances that are exclusive to said thread(s). This design moves the area in which inter-thread communication is required from the API instances to the dispatch layer. An advantage over the direct access architecture is the massively reduced degree of contention on the API instances, as only ever a single thread will access the SQ or CQ, respectively, which is especially crucial for the CQ, due to the polling nature of the access. Concurrent polling for completions of requests from multiple users thus does not reduce I/O throughput, as users poll on the objects they received for this purpose and not directly on the CQ. Furthermore, while significant contention could occur during submission to the dispatch layer, the design of the dispatch layer can be adapted to address this if required, e.g. by increasing the number of data structures used to buffer and distribute the requests.

While resources allocated for I/O and non-I/O can be scaled independently, similar to the direct access architecture, but decoupled further as the CPU time for I/O tasks is provided by the dedicated threads of the I/O instances, this allocation is still static. Furthermore, as the threads of I/O Instance are dedicated threads, it is possible to utilize these as I/O threads. Implementations could also as an additional optional mechanism choose to provide the option to utilize inline execution via callbacks to prioritize latency-sensitive tasks. While this architecture provides many benefits the potentially significant increase in inter-thread communication compared to the shared-nothing approach remains its most significant trade-off. The extent of this overhead depends on the granularity of submitted tasks and on the many implementation details of the dispatch layer and I/O instances. These design decisions range from the used data structures and synchronization primitives in the dispatch layer over the used scheduling scheme to the mechanism used to manage the metadata of in flight requests.

A system that can be categorized in this architecture is Go's scheduler for *goroutines*, which are, like KSEs, stackful, thus reducing the performance issues due to nesting. Go's usability in the context of high performance async I/O is, however, debatable, due to its managed nature, most notably, but not limited to its garbage collector. An academic system that similarly should be classified under this architecture is Leanstore [11], which utilizes a single queue as the dispatch layer and callback task partitioning. Merzljak et al. [23] base their work on Leanstore, thus sharing its architecture, but utilize C++ coroutines instead.

### 4.2.4 Dynamic I/O Threadpool

Notably, Haas et al. discussed the impact of the SQ poll threads in their design and concluded that it would be more beneficial to turn off the feature entirely [11]. We argue that given the extent to which both SQ polling can increase per-

formance [4, 6, 20] further exploration of the management of these kernel threads is warranted. If the design is adapted such that the amount of I/O resources, including these kernel threads, are accurately scaled to what is actually required, both statically to the hardware capabilities and dynamically to the current I/O load, the performance gains of SQ polling could be leveraged while retaining good or even better CPU efficiency. To this end two issues of the previous architecture would need to be addressed. The first is to utilize full task partitioning as the primary partitioning scheme over callback partitioning or coroutines. Besides the established trade-offs, this allows the amount of resources allocated to I/O to be accurately scaled, e.g., statically to conform to hardware capabilities, as I/O instances now *only* perform I/O-related work. If, as in § 4.2.3, hybrid execution is allowed, assumptions on the average I/O to non-I/O work ratio have to be made to estimate the number of I/O instances that would not bottleneck the hardware, with any overshoot in the estimate resulting in reduced CPU efficiency.

Similarly dynamically scaling the amount of active I/O instances based on the current I/O load is now efficiently possible, as with IO threads reducing the amount of active I/O instances does not proportionally reduce the amount of resources available for non-I/O related tasks. Deactivating instances can be achieved by adjusting the distribution of requests in the distribution layer to skip the I/O instances in question. If an I/O instance is consequently starved of requests for long enough, the corresponding kernel thread goes to sleep, and the thread(-pair) of the I/O instance should be implemented to then also go to sleep. Alternatively the CPU time of the user space threads of an I/O instance could also be provided by the users via a function call, rather than a dedicated thread. In this case, rather than going to sleep, the thread or thread-pair could simply return. This approach distributes requests to fewer I/O instances in times of low load and releases the high upkeep cost resources of polling threads of currently unused I/O instances.

## 5 Summary

While contemporary async APIs, such as `io_uring`, can significantly improve performance, their adoption both in commercial and academic systems has been limited, likely due to the complexities inherent in achieving adequate performance when using them, which have been laid out in this paper. We discussed the most notable architectures that may be used for high performance async I/O. The shared nothing approach, which, although most efficient, is limited in applicability, the static I/O thread pool, which struggles with scaling and CPU efficiency, if otherwise beneficial features such as SQ polling are used and the dynamic I/O thread pool that attempts to address this issue by actively managing the kernel threads introduced by SQ polling. We hope that our analysis of currently employed asynchronous I/O architectures will

spark further discussion and lead to a better adoption and more convenient integration of asynchronous I/O APIs into I/O-intensive applications, such as database systems and I/O libraries.

## References

- [1] `io_uring` - win32 apps.
- [2] Specifications - NVM express.
- [3] AMD. 4th gen amd epyc processor architecture. White paper, AMD, 5 2024.
- [4] AXBOE, J. Efficient IO with `io_uring`, 10 2019. [Online; accessed 20-October-2023].
- [5] CORRETTÉ, S. I. `libuv`, 2024.
- [6] DIDONA, D., PFEFFERLE, J., IOANNOU, N., METZLER, B., AND TRIVEDI, A. Understanding modern storage APIs: a systematic study of `libaio`, `SPDK`, and `io_uring`. In *Proceedings of the 15th ACM International Conference on Systems and Storage* (New York, NY, USA, June 2022), SYSTOR '22, pp. 120–127.
- [7] ET AL., A. K. `Zig`, 2024.
- [8] ET AL., C. L. `tokio-uring`, 2024.
- [9] GREEF, J. D. `Tigerbeetle`, 2024.
- [10] GU, J., WU, X., LI, W., LIU, N., MI, Z., XIA, Y., AND CHEN, H. Harmonizing performance and isolation in microkernels with efficient intra-kernel isolation and communication. In *USENIX ATC'20* (2020).
- [11] HAAS, G., AND LEIS, V. What Modern NVMe Storage Can Do, and How to Exploit It: High-Performance I/O for High-Performance Storage Engines. *Proc. VLDB Endow.* 16, 9 (2023), 2090–2102.
- [12] HE, Y., LU, J., AND WANG, T. `CoroBase`: coroutine-oriented main-memory database engine.
- [13] KIOXIA. `Cd8p-r` product brief. Product brief, Kioxia, 2024.
- [14] KIVITY, A. `Seastar`, 2024.
- [15] KOCBERBER, O., FALSAFI, B., AND GROT, B. Asynchronous memory access chaining. *Proc. VLDB Endow.* 9, 4 (12 2015).
- [16] LEIS, V., HAUBENSCHILD, M., KEMPER, A., AND NEUMANN, T. `LeanStore`: In-Memory Data Management beyond Main Memory. In *ICDE'18* (2018).
- [17] LUND, S. A. F., BONNET, P., JENSEN, K. B. A., AND GONZALEZ, J. I/O interface independence with `xnvme`. In *SYSTOR '22* (2022).
- [18] MI, Z., LI, D., YANG, Z., WANG, X., AND CHEN, H. `SkyBridge`: Fast and Secure Inter-Process Communication for Microkernels. In *EuroSys '19* (2019).
- [19] PHISON. Product brochure `ps5026-e26`. Product brief, Phison, 2024.
- [20] REN, Z., AND TRIVEDI, A. Performance Characterization of Modern Storage Stacks: `POSIX I/O`, `Libaio`, `SPDK`, and `Io_uring`. In *CHEOPS '23* (2023).
- [21] SEMICONDUCTOR, S. `Pm1743 ssd` whitepaper. White paper, Samsung Semiconductor, 2021.
- [22] SOARES, L., AND STUMM, M. `FlexSC`: Flexible system call scheduling with Exception-Less system calls. In *OSDI'10* (2010).
- [23] VON MERZLJAK, L., FENT, P., NEUMANN, T., AND GICEVA, J. What are you waiting for? use coroutines for asynchronous I/O to hide I/O latencies and maximize the read bandwidth! In *ADMS'22* (2022).
- [24] YANG, Z., HARRIS, J. R., WALKER, B., VERKAMP, D., LIU, C., CHANG, C., CAO, G., STERN, J., VERMA, V., AND PAUL, L. E. `SPDK`: A Development Kit to Build High Performance Storage Applications. In *CloudCom'17* (2017).
- [25] ZHOU, Z., BI, Y., WAN, J., ZHOU, Y., AND LI, Z. `Userspace Bypass`: Accelerating Syscall-intensive Applications. In *OSDI'23* (2023).