

Comparing User, Kernel, & Kernel-Bypass HTTP Services

Abstract

Low-latency networking is essential in modern computer services. Distributed networking and online gaming represent use cases that require data to be sent, processed, and received with very low overhead to achieve maximum performance. Our experiments seek to analyze the empirical performance characteristics of servicing network requests at different boundary levels of the operating system, in particular servicing packets with user-level processes, kernel-level modules, and direct kernel-bypass NIC communication in the data plane.

Background and Problem

Network communication is ubiquitous in modern computing. For standard user-level applications in a monolithic kernel, network communication requires the process to communicate with the OS abstraction layer, which then interacts with the network interface card (NIC) driver, which then sends instructions to the NIC hardware to send a packet. The reverse process happens to receive packets. These layers of indirection increase overall latency to process and respond to packets, adversely affecting performance in critical low-latency high-communication applications such as grid computing, competitive gaming, or high-frequency trading.

Many system designs have been proposed, implemented, and tested to reduce latency and improve performance while still maintaining an abstraction layer to share the network resource and interface among multiple processes. JITSU treats incoming network packets as signals for servers to rapidly spin up unikernels to process the packet and send the proper response, but faces higher latency due to kernel boot delay. Arrakis sets up access for a process to have direct dataflow to the NIC for a time quantum before interposing control to another process, allowing fast response for the process with exclusive packet access, but requiring other processes to wait their turn. Exokernel uses injected kernel code to process packets based on a filter without having to switch on the user boundary, prioritizing speed over security. SnapOS uses a user-level library to bypass the kernel with a separate channel to the hardware NIC, ensuring the path to network access is short, but requiring explicit support from the NIC driver through SR-IOV and explicit cooperation of user processes.

Our goal is to investigate the performance of the underlying methods used in these systems in terms of throughput and latency in responding to HTTP requests. We look at the performance of some modern userspace web servers, Apache and Flaskr, as a baseline. We then compare this to a simple web server daemon implemented within the Linux kernel as a kernel module (khttpd) and an application that interfaces directly to the NIC via SR-IOV using the Data Plane Development Kit (DPDK) (with a lot of help from tinyhttpd-lwip-dpdk).

Hypothesis

Before starting, we hypothesized that in terms of latency, Flaskr will perform worst, followed by the kernel module server, with the DPDK server performing the fastest, as each layer removes a level of abstraction overhead from userspace to kernel, then from kernel to hardware interface. While Apache is theoretically also slower, we reserve the possibility that it could perform better than anticipated and potentially better than a kernel module due to its widespread usage and continuous optimization under the worse-is-better principle.

The Journey and Setbacks

For the kernel module, we first found existing implementations of khttpd as experimental code ([Github](#)). This did not have the ability to read files or serve concurrent requests with multiple threads, so we found a forked version that implemented that functionality (mostly) to use for testing ([Github](#)). We first attempted to compile and run the module on WSL, but WSL's kernel does not come with Linux headers that are needed for the kernel module interface. Next, we tried to directly inject the installed module into the kernel used in lab0 so it starts on boot. That also did not work (the server did not start on boot), but using the `insmod` install module command did, enabling curl commands to be sent on localhost and receiving directory listings. However, the module crashes whenever a file is attempted to be read, and the kernel messages that resulted were difficult to parse due to a lack of error messages.

```

Machine View
168.0689201 CR2: 0000000000000000 CR3: 000000002640a005 CR4: 000000000360ee0
168.0701091 Call Trace:
168.0706061 ? __die+0x81/0xc3
168.0711331 ? no_context+0x195/0x380
168.0717421 ? __do_page_fault+0xb2/0x4f0
168.0724121 ? do_filp_open+0xa7/0x100
168.0730481 ? async_page_fault+0x1e/0x30
168.0737671 ? filp_open+0x3f/0x50
168.0743481 ? rw_verify_area+0x10/0xb0
168.0750121 ufs_read+0x6b/0x140
168.0758751 kernel_read+0x2c/0x40
168.0764771 handle_directory+0x236/0x460 [khttpd]
168.0772371 http_parser_callback_message_complete+0x17/0x30 [khttpd]
168.0782231 http_parser_execute+0x3272/0x3450 [khttpd]
168.0790341 http_server_worker+0x162/0x1d0 [khttpd]
168.0798311 ? http_parser_callback_request_url+0x20/0x20 [khttpd]
168.0813261 ? http_server_send+0xb0/0xb0 [khttpd]
168.0820801 ? http_parser_set_max_header_size+0x10/0x10 [khttpd]
168.0833731 ? tracedir+0x50/0x50 [khttpd]
168.0840751 ? http_parser_callback_header_field+0x10/0x10 [khttpd]
168.0850241 ? http_server_worker+0x1d0/0x1d0 [khttpd]
168.0881501 ? handle_directory+0x460/0x460 [khttpd]
168.0892521 ? tracedir.part.2+0x100/0x100 [khttpd]
168.0900681 process_one_work+0x1a7/0x3a0
168.0907421 worker_thread+0x30/0x390
168.0913611 ? create_worker+0x1a0/0x1a0
168.0920051 kthread+0x112/0x130
168.0945781 ? kthread_bind+0x30/0x30
168.0955051 ret_from_fork+0x1f/0x40
168.0963631 Modules linked in: khttpd(OE) nls_ascii nls_cp437 ufat fat ppdev kum_intel bochs_drm ttn kum irqbypass drm_kms_helper
parport_pc sg serio_raw parport qemu_fw_cfg button ip_tables x_tables autofs4 ext4 crc16 mbcache jbd2 crc32c_generic fs/crypto ecb sd_mod sr_mod cdrom crc32c_intel ata_generic ata_piix
libata aesi intel scsi_mod psmouse aes_x86_64 crypto_sind cryptd glue_helper floppy e1000 i2c_piix4
168.1050401 CR2: 0000000000000000
168.1059011 ---[ end trace fea778adb4f6f067 ]---
168.1073041 RIP: 0010:rw_verify_area+0x10/0xb0
168.1083831 Code: e8 75 36 00 00 eb a1 48 c7 c3 f7 ff ff eb 98 e8 95 2b e2 ff 0f 1f 44 00 00 0f 1f 44 00 00 48 8b 46 20 48
85 c9 78 36 55 53 <8> 8b 12 48 85 d2 78 32 48 89 43 48 01 cb 78 7c 48 83 b8 70 01 00
168.1147871 RSP: 0018:ffffb54100353970 EFLAGS: 00010202
168.1158861 RAX: ffff98e3780200e8 RBX: 0000000000000010 RCX: 0000000000000010
168.1173041 RBX: 0000000000000000 RSI: ffff98e378bf3400 RDI: 0000000000000000
168.1187131 RBP: ffffffff98e378bf3400 R08: ffffffff98e378bf3400 R09: ffffffff98e378bf3400
168.1201021 R10: 0000000000000000 R11: ffff98e335ef7b40 R12: ffff98e378bf3400
168.1214741 R13: 0000000000000001 R14: ffff98e3775bba20 R15: 0000000000000000
168.1228381 FS: 0000000000000000 (0000) GS: ffff98e37ab00000 (0000) kn1GS:0000000000000000
168.1243601 CS: 0010 DS: 0000 ES: 0000 CR0: 0000000000005003
168.1256521 CR2: 0000000000000000 CR3: 000000002640a005 CR4: 000000000360ee0

```

Image 1: Kernel messages from khttpd crashing on file GET request.

We eventually found the root cause that separated the behavior of reading directories from reading files was the `kernel_read` call, which requires a change to pass a pointer to 0 rather than a 0 (which is incorrectly interpreted as a null pointer). We have not pushed this change upstream (so good luck to everyone else figuring it out). With a functional server on the VM, we worked to port it to the physical machine serving as the benchmarking environment. This presented a new challenge in the form of security: because kernel modules have broad unrestricted access, installing a module was blocked even with superuser privileges.

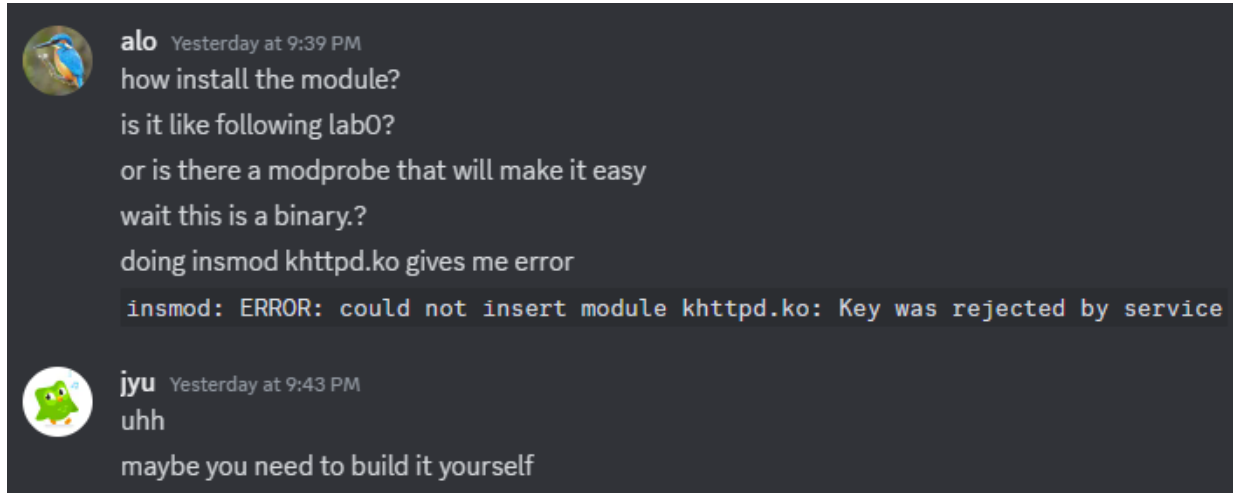


Image 2: Attempting to resolve the `insmod` error. Building it again did not help.

After further research, we found the error was due to secure boot. We also found along the way that the VM Linux and physical machine Linux versions differed, and the rebuild process required additional patching to handle gcc version differences that resulted. After temporarily disabling the option in the BIOS, the module could then be loaded into the kernel while running and we could benchmark it. A limitation of this kernel server is that the entire file being read is first copied into kernel memory and then sent via HTTP, so there is a maximum limit to the size of the file in this prototype that commercial web servers address.

For the implementation of an HTTP server with DPDK as the network backend, the first step was to set up the system so that the NIC is able to communicate with the internal system. In order to do this, the NIC must be placed in its own IOMMU group in order to bind the ports to VFIO for DPDK usage. However, it was discovered that the placement of the Intel NIC was actually in an IOMMU group that contained multiple other PCI devices, including the storage device that contained the currently running OS. Deciding that that was unwise (and finding other ways of breaking IOMMU groups questionable), we eventually decided to take apart the machine and strategically move PCI devices around to different points on the motherboard until everything was arranged in a way such that the NIC was in its own IOMMU group.

```

vespaston@vespaston-MS-7C56:~/Downloads/dpdk-23.11/examples/helloworld$ sudo ./build/helloworld
EAL: Detected CPU lcores: 12
EAL: Detected NUMA nodes: 1
EAL: Detected shared linkage of DPDK
EAL: Multi-process socket /var/run/dpdk/rte/mp_socket
EAL: Selected IOVA mode 'VA'
EAL: VFIO support initialized
EAL: 0000:04:00.0 VFIO group is not viable! Not all devices in IOMMU group bound to VFIO or unbound
EAL: Requested device 0000:04:00.0 cannot be used
EAL: 0000:04:00.1 VFIO group is not viable! Not all devices in IOMMU group bound to VFIO or unbound
EAL: Requested device 0000:04:00.1 cannot be used
TELEMETRY: No legacy callbacks, legacy socket not created
hello from core 1
hello from core 2
hello from core 3
hello from core 4
hello from core 5
hello from core 6
hello from core 7
hello from core 8
hello from core 9
hello from core 10
hello from core 11
hello from core 0
vespaston@vespaston-MS-7C56:~/Downloads/dpdk-23.11/examples/helloworld$ █

```

Image 3: DPDK complaining about IOMMU groups

After that, we were able to not-easily set up a DPDK script that forwarded packets received from SR-IOV NIC to an internal script that polled continuously for packets. While this seemed to work great in the sense that packets could be sent back and forth through this bridge, and the NIC was forwarding packets to the application, we could not intentionally ping/send packets to said application through the NIC (internal DPDK packet generators could communicate but not external). This probably occurred due to the lack of backend or in-depth network issues that remain outside our knowledge. It was then at this point that we realized that DPDK does not provide a TCP/IP stack, requiring the user to build it themselves. As this would be insurmountable for us to complete in the short time frame we had, we turned to external libraries that promised to help. One that we chose was SeaStar ([Github](#)), which promised to be a high-performance networking library. After spending much time trying to get this to work, it still throws countless errors trying to compile. We remain very confused as to how others got it working as even the mere import of their libraries causes gcc to complain about a myriad of issues and missing variables. Finally, we settled on another project tinyhttpd ([Github](#)) that combines both DPDK with lwIP, a lightweight TCP/IP stack, to build a simple web server. As we were still managing the entire TCP/IP stack, dealing with HTTP splitting into multiple TCP packets was out of scope, thereby limiting it to a single TCP packet and limiting the size of a possible HTTP response. DPDK still remains very confusing to us and perhaps was not the best thing to attempt in this project.

Limitations: Limited HTTP size therefore our HTTP server testing was restricted to small files.

Environment

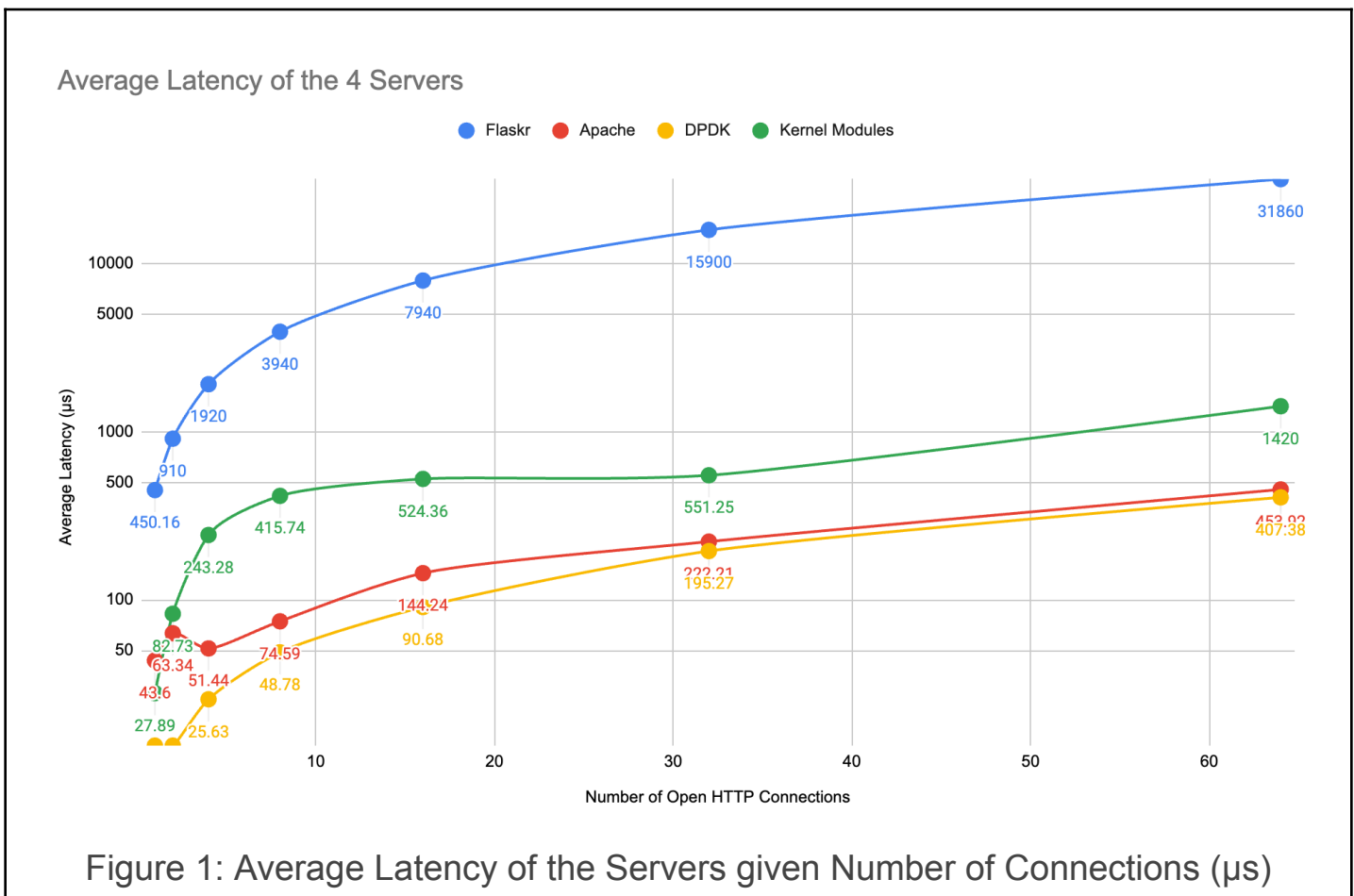
These experiments were all performed on a PC that was running Ubuntu-22.04. The PC was running an AMD Ryzen 5 5600X 6-Core Processor that has support for AMD-V virtualization for an x86_64 architecture and 12 CPUs. In addition, it was running on top of a B550-A Pro motherboard that had support for SR-IOV PCIe. The two network cards within the system were

a Realtek RTL8111 PCI Express Gigabit Ethernet Controller and an NIC with an Intel 82576 card (that supports SR-IOV as well as is supported by DPDK).

Experimental Design

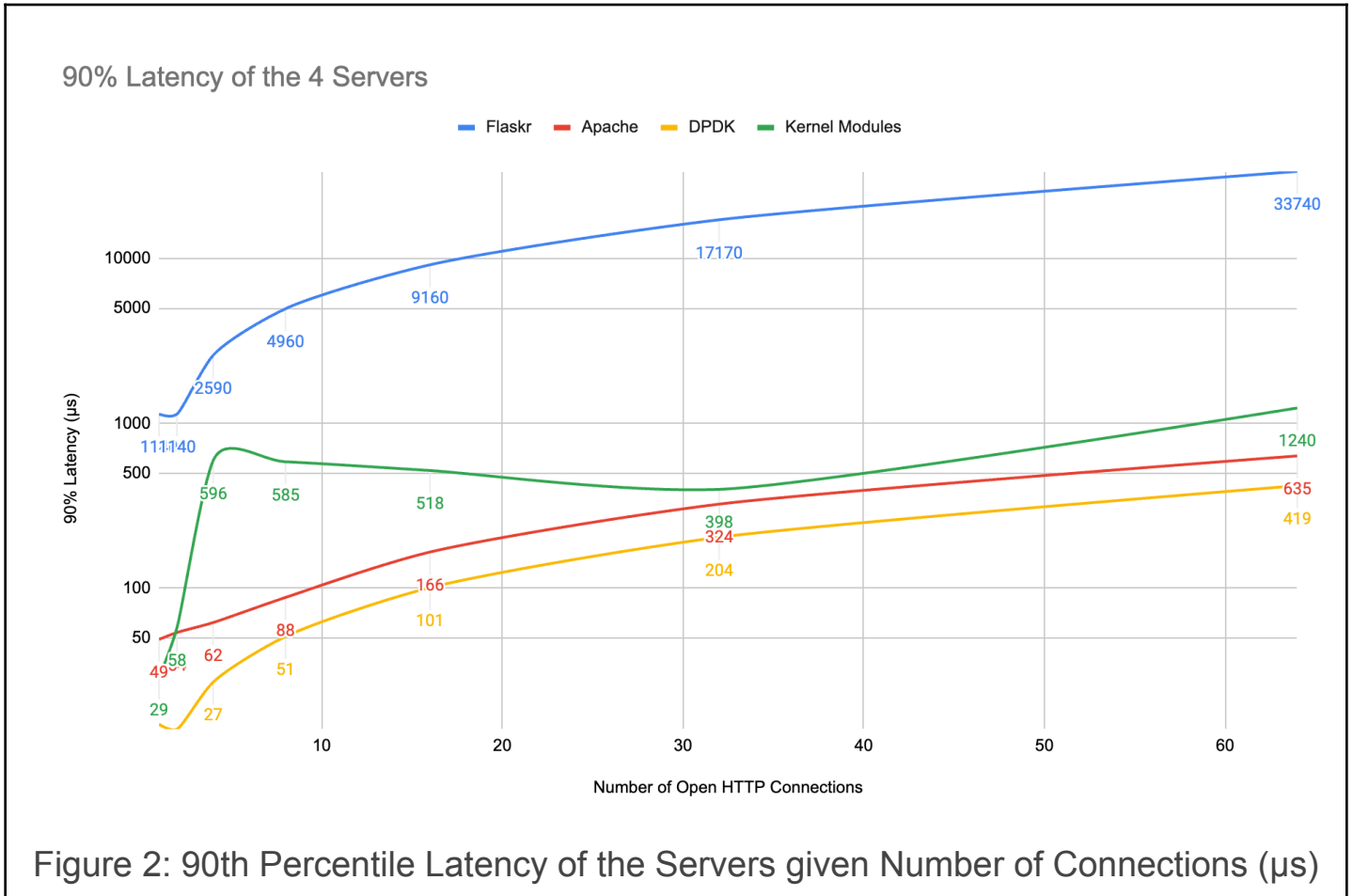
Due to the limitations discussed above, we were only able to test for latency on small files. Stemming from this, we decided to use the HTTP benchmarking tool [wrk](#) to test the latency of retrieval of a certain file. We utilize *wrk* to run a benchmark that runs for 10 seconds and keeps a varied number of HTTP connections open to the server. The number of these connections go from 1 to 64 in powers of 2 and we measure the average latency, 90% latency, and request throughput of the server.

Results

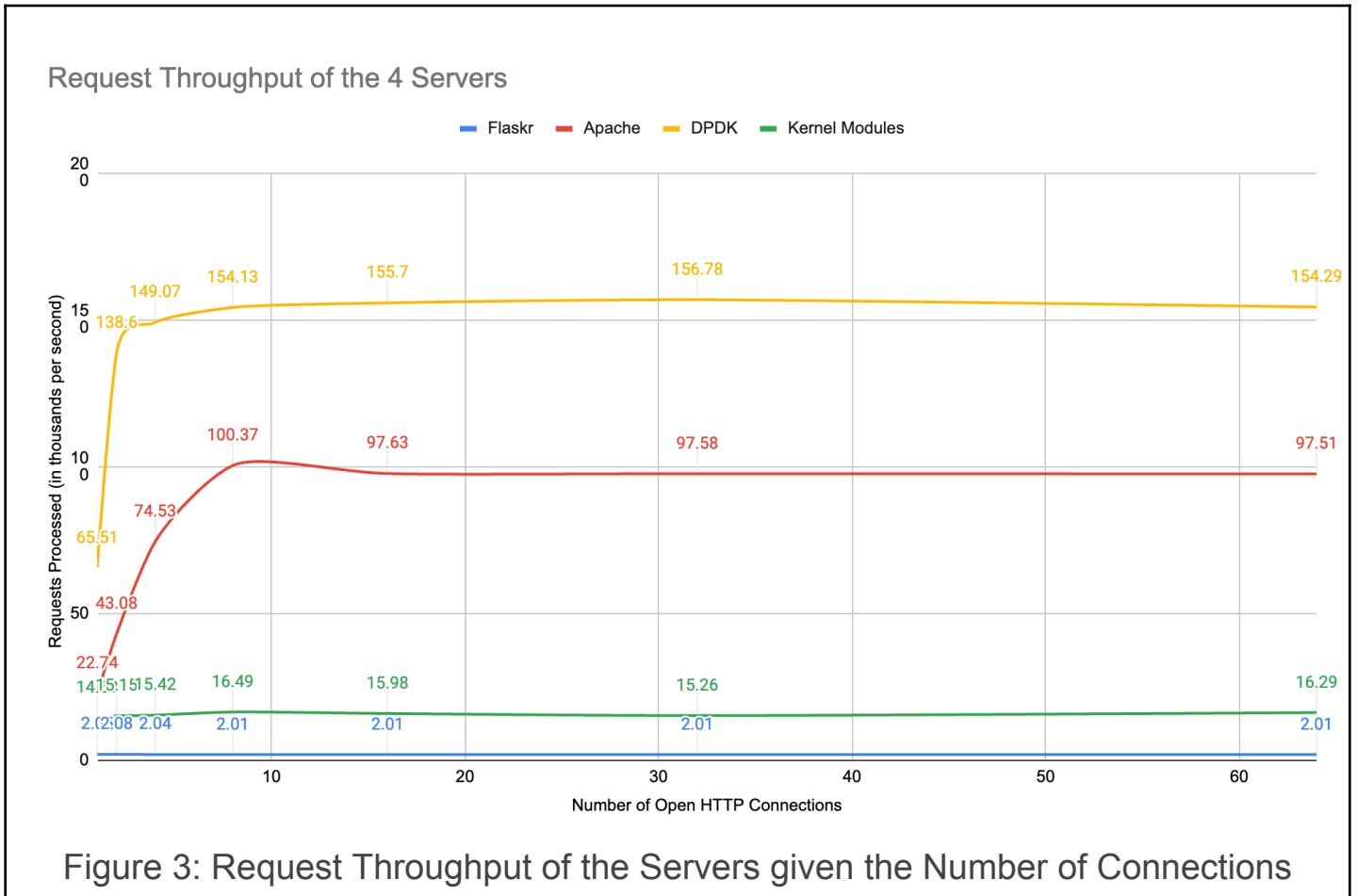


When comparing average latency, the results (conveniently) support our hypotheses, with Flaskr having the worst latency followed by the kernel module and the DPDK service. As anticipated, due to Apache's long usage and continuous development, its performance is very good, better than the kernel server we used, and gets close to matching the DPDK setup at high connection

amounts. This suggests that Apache is optimized to scale to high numbers of requests at once and handle load-stress scenarios.



The results of the 90th percentile latency response we received from the servers matches up very well with our previous results of the overall average latency. If a server experiences a lower average latency, so will its 90th percentile latency response. However, there were still several interesting things to note. Comparing Apache and DPDK, the difference between DPDK average latency and 90th percentile latency is quite low, suggesting packets experience a very uniform treatment through the pipeline, which makes sense as the interrupts that normally occur are mostly removed, whether that be an interrupt from new packet received or syscalls required to interact with network stack. Comparatively, Apache's is much higher, given the reasons above. Another thing to note is that the Kernel Modules implementation faced absolutely horrible 99% that hit up to 30 times worse latency compared to the average. Why this is an issue specific to the kernel implementation compared to userspace and hardware implementations remains a mystery.



The results here matched what we expected from the other measures of average latency and 90% latency. The implementations that performed well there perform well here. One interesting thing to note is how quickly it took to saturate the system given the open HTTP connections for DPDK versus Apache. DPDK very quickly got saturated at around 3 connections while Apache lasted up until 10 connections before reaching saturation. This again suggests that Apache has been optimized for high requests, given its wide adoption in server environments.