



Energy Implications of IO Interface Design Choices

Sidharth Sundar, William Simpson, Jacob Higdon,
Caeden Whitaker, Bryan Harris, Nihat Altiparmak

Computer Systems Lab

Computer Science & Engineering Department

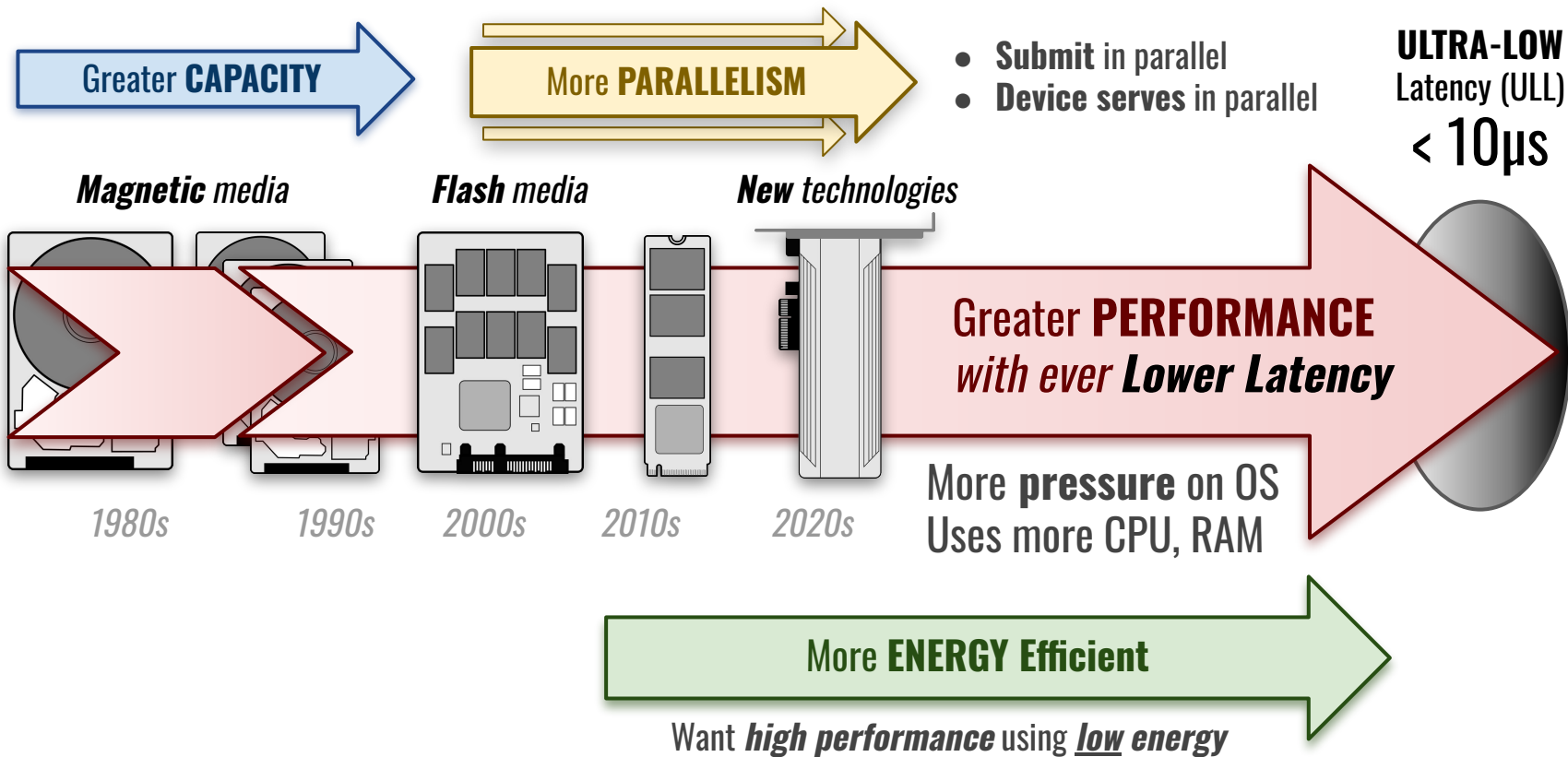
UNIVERSITY OF LOUISVILLE[®]

HotStorage '23 (Boston), July 9, 2023

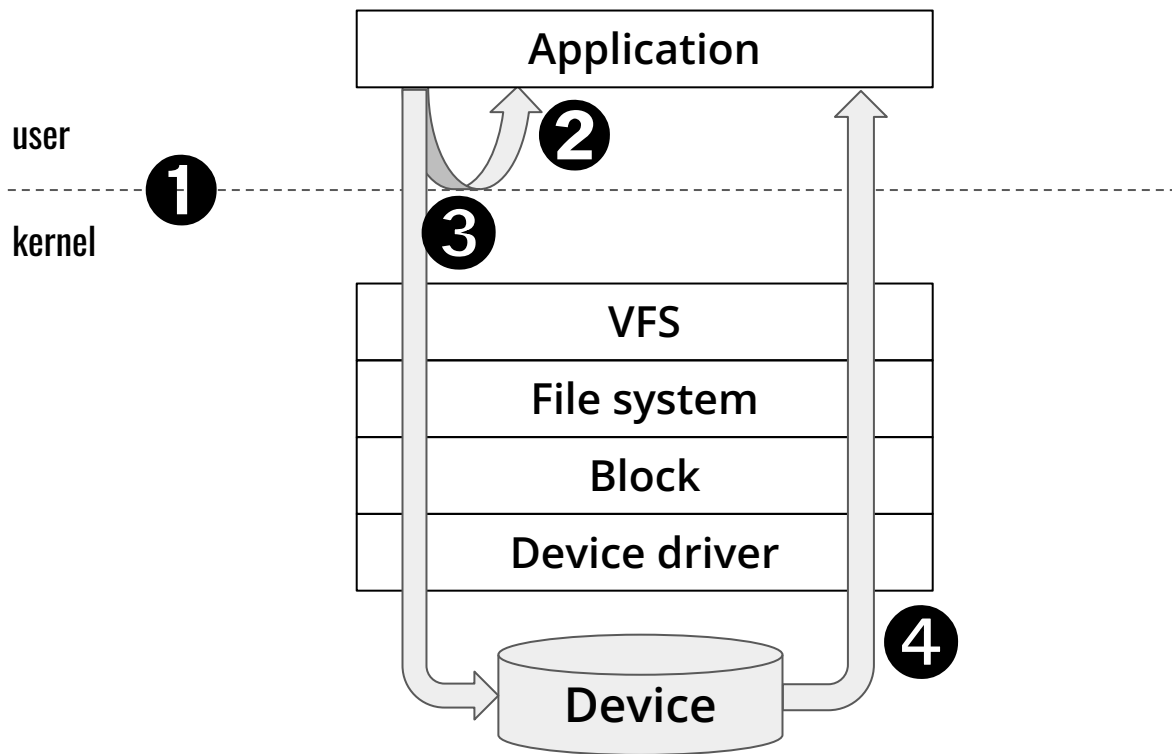
Outline

- Motivation
- IO Interface design choices
- Experimental Results
 - Latency Impact
 - Energy Efficiency
- Conclusion

Trends in data storage



IO Interfaces



IO Interface Design Choices

Existing Linux APIs

posix-sio — Based on traditional POSIX synchronous *read(2)*, *pread(2)*, etc.

posix-aio — POSIX asynchronous, implemented in library on top of *posix-sio*.

libaio — Linux native asynchronous IO library

io_uring — Recent Linux design for high performance

spdk — Intel framework based on kernel bypass

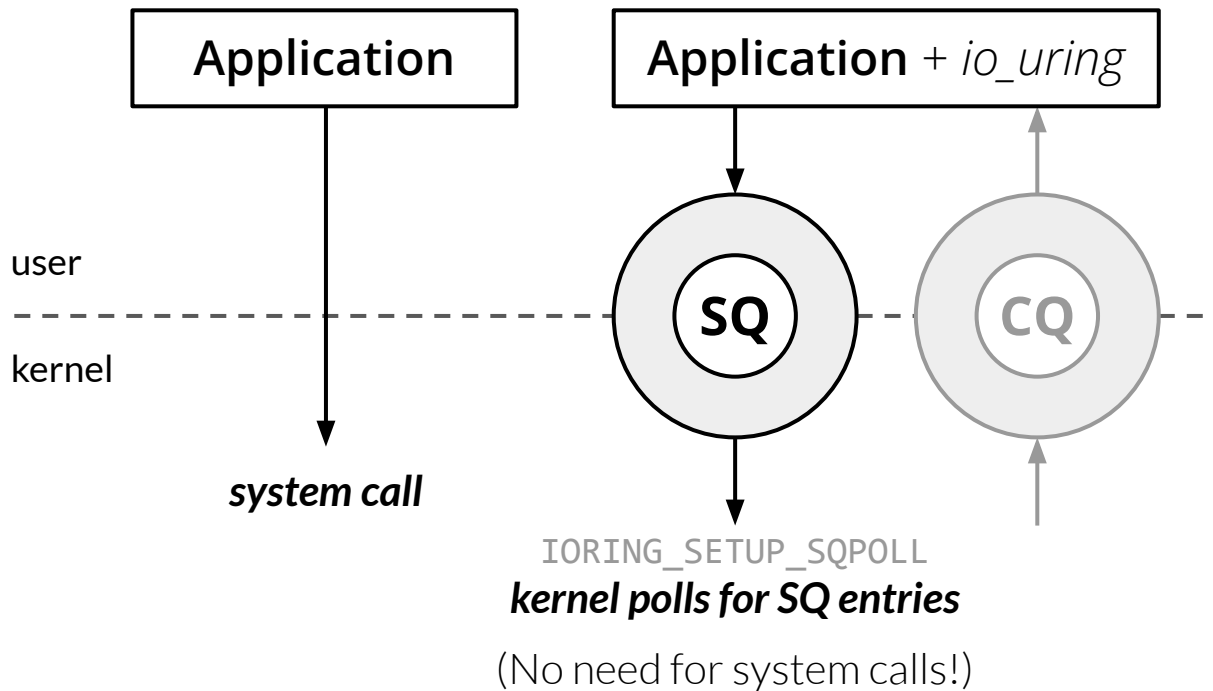
❶ Execution in *kernel* or *user* space

API	kernel	user
<i>posix-sio</i>	Implemented in kernel	
<i>posix-aio</i>		Implemented in C lib. on top of <i>posix-sio</i>
<i>libaio</i>		
<i>io_uring</i>		
<i>spdk</i>		Uses kernel bypass drivers

② Synchronous or Asynchronous behavior

API	Synchronous	Asynchronous
<i>posix-sio</i>	synch. only	
<i>posix-aio</i>		asynch. only
<i>libaio</i>		asynch. only
<i>io_uring</i>	both supported	
<i>spdk</i>		asynch. only

③ **Submission** using *system call* or *polling*



③ **Submission** using *system call* or *polling*

API	System call	Submission polling
<i>posix-sio</i>	System call based	
<i>posix-aio</i>		
<i>libaio</i>		
<i>io_uring</i>		<i>io_uring</i> feature
<i>spdk</i>	kernel bypass	

④ Completion using *interrupt* or *polling*

API	Interrupt	Poll
<i>posix-sio</i>	Interrupts	Optional polling
<i>posix-aio</i>		
<i>libaio</i>		
<i>io_uring</i>		Optional polling
<i>spdk</i>		Polling required

Interface design choices

API	Execution		Behavior		Submission		Completion	
	kernel	user	sync	async	syscall	poll	int	cpoll
<i>posix-sio</i>	●	○	●	○	●	○	●	●
<i>posix-aio</i>	●	●	○	●	●	○	●	○
<i>libaio</i>	●	○	○	●	●	○	●	○
<i>io_uring</i>	●	○	●	●	●	●	●	●
<i>spdk</i>	○	●	○	●	○	○	○	●

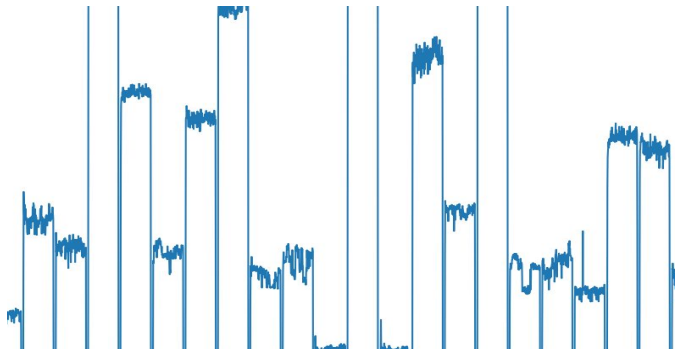
○ unsupported, ● supported

Experimental Results

Experimental Setup



Image from Onset Computer Corporation



Power Measurement

Onset HOB0 plug meter logs power, current, etc., every second, for the entire system.

Workloads

fio ("Flexible IO tester")

– direct IO

- *xfs* file system
- "none" IO scheduler

Latency Impact

Latency Impact

Interface	Subm.	Comp.	read (μ s)		write (μ s)		avg (%) slower than <i>spdk</i>
			50th	99th	50th	99th	
<i>spdk</i>	—	cpoll	6.23	7.47	8.44	10.40	—
<i>posix-sio</i>	syscall	cpoll	8.69	9.06	11.11	12.24	27.52%
<i>io_uring</i>	syscall	cpoll	8.65	9.03	14.50	15.93	46.18%
<i>io_uring</i>	spoll	cpoll	7.89	9.05	12.38	22.33	52.30%
<i>posix-sio</i>	syscall	int	11.68	12.36	14.08	15.64	67.54%
<i>io_uring</i>	spoll	int	8.72	11.15	13.34	25.32	72.69%
<i>libaio</i>	syscall	int	12.35	13.06	16.33	21.87	94.21%
<i>io_uring</i>	syscall	int	12.73	13.45	16.92	23.61	102.97%
<i>posix-aio</i>	syscall	int	19.10	28.89	21.60	34.46	220.15%

Energy Impact

Energy efficiency metric

$$\text{Energy efficiency} = \frac{\text{Performance}}{\text{Power}} = \frac{\text{IOPS}}{\text{Watts (J/s)}} = \text{IO/J}$$

Energy Impact Experiments

- 4KB, 16KB, and 128KB requests
- 100% random read vs 100% random write
- Scaling # of overall requests issued in parallel to device using 2 metrics:

1. Single Thread

- a. Scaling # of outstanding requests (iodepth) on single thread
- b. Single thread issuing IOs
- c. Only considers asynchronous interfaces

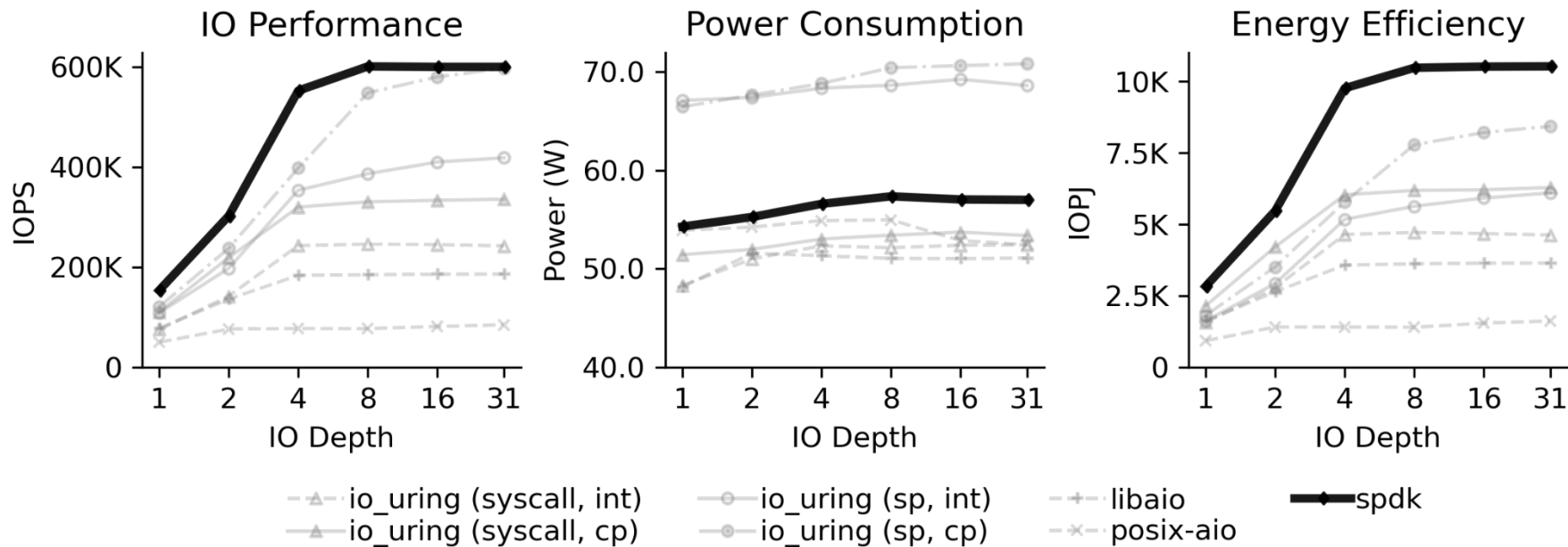
2. Multi-threaded

- a. Scaling # of threads issuing IOs (# of jobs)
- b. One request per thread

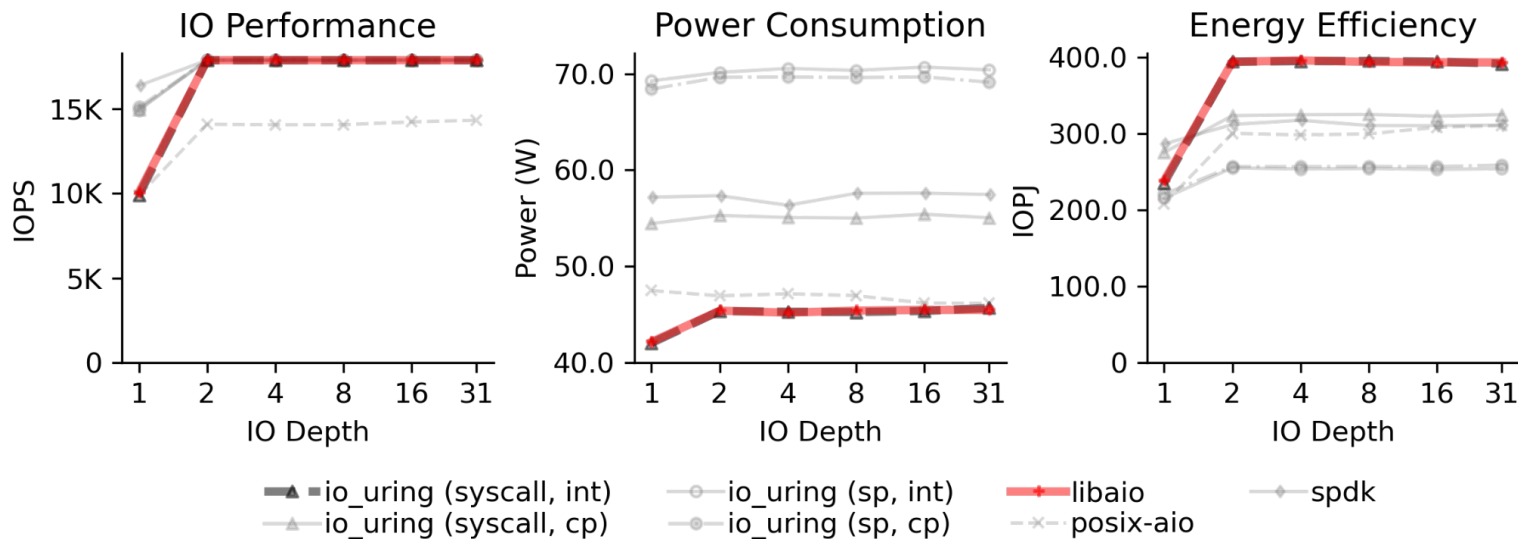
Energy Impact

Single Thread

Kernel bypassing is most energy efficient for single thread small requests



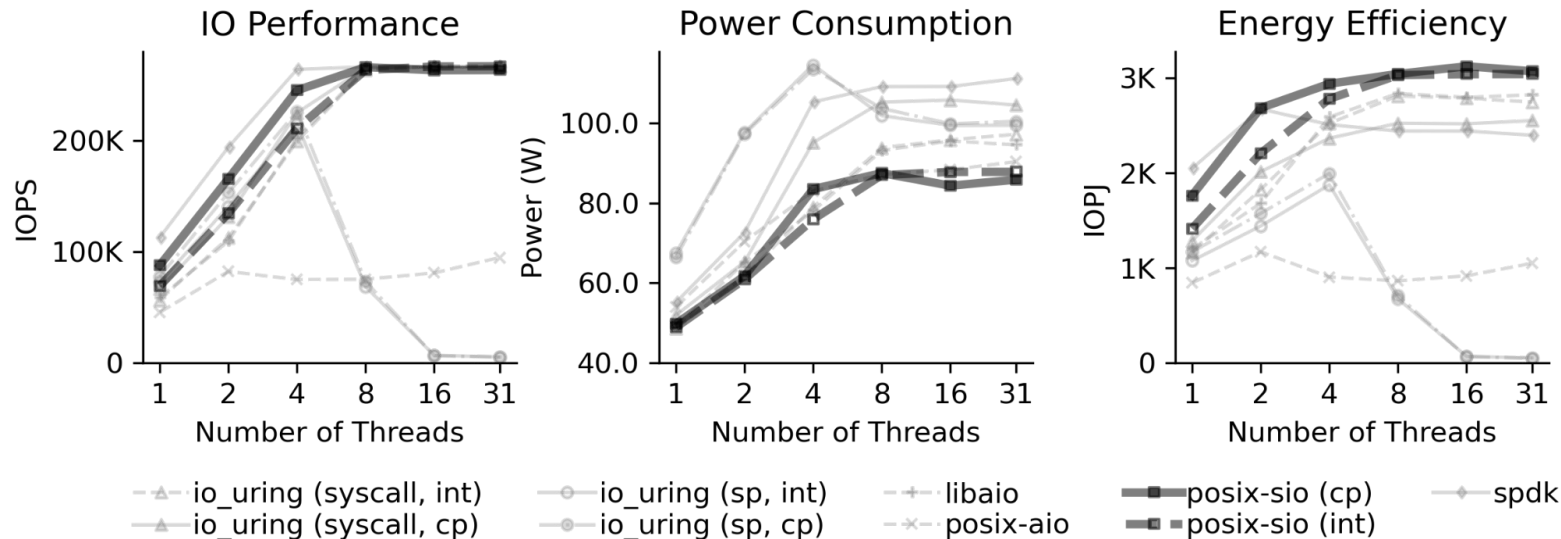
System call and interrupt based kernel space implementations are most efficient for large requests



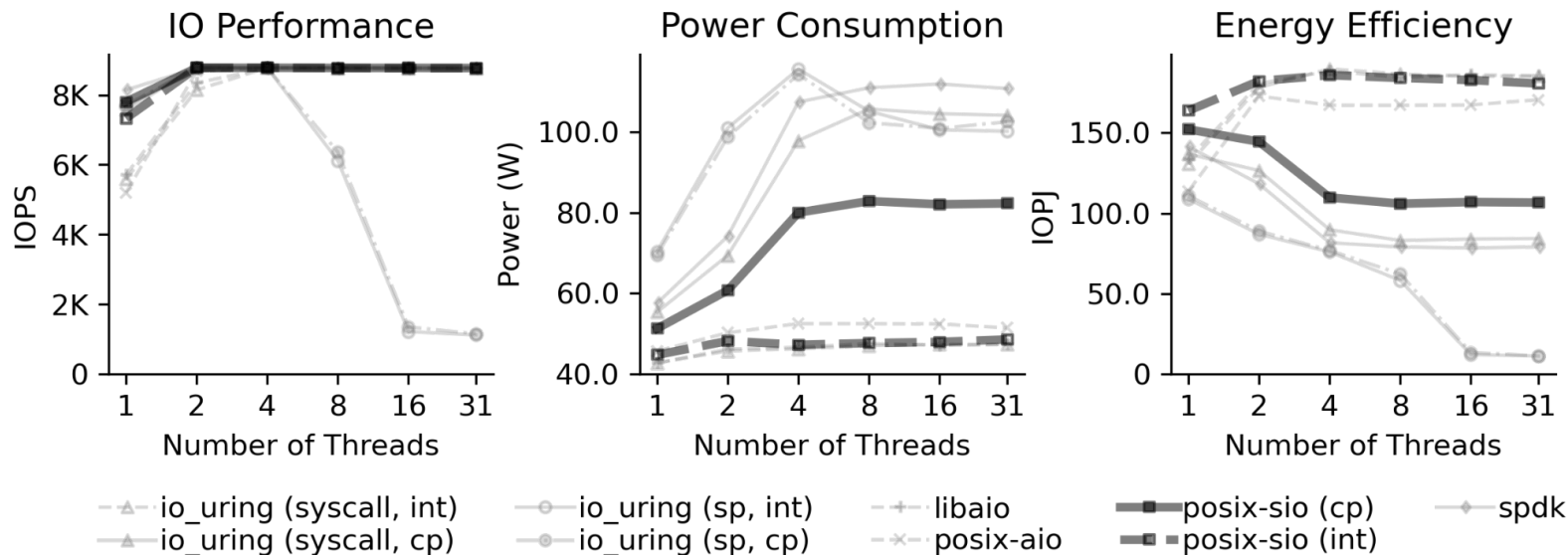
128 KB Single Thread Random Reads

Energy Impact Multithreaded

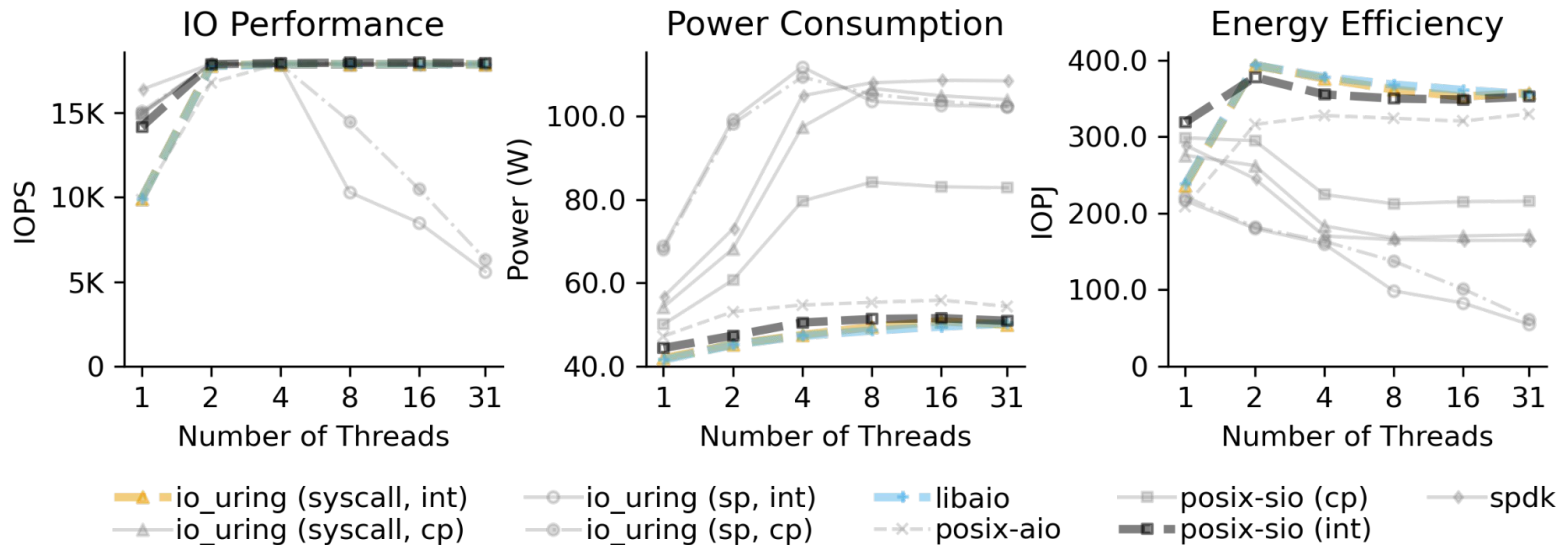
For small requests, *posix-sio* with polling based completion is most energy efficient



For large requests, *posix-sio* with interrupt based completion is most energy efficient



Interrupts are crucial for energy efficiency when the request size gets larger



128 KB Multithreaded Random Reads

Conclusion

① Kernel *VS* Userspace

Only bypass kernel if kernel functionalities (such as interrupts) are unnecessary

② Synchronous *VS* Asynchronous

Synchronous (posix-sio) tends to be more energy efficient when synchronous is usable

③ System Call *VS* Submission Polling

Submission polling typically costs too much power to justify

④ Interrupt *VS* Completion Polling

- **Polling** is more energy efficient for **smaller** IOs
- **Interrupt** is more energy efficient for **larger** IOs

Questions?