Arrakis: The Operating System is the Control Plane

Simon Peter, Jialin Li, Irene Zhang, Dan R. K. Ports, Doug Woos, Arvind Krishnamurthy, Thomas Anderson, Timothy Roscoe*

University of Washington, *ETH Zurich

Symposium on Operating Systems Design and Implementation (OSDI), October 2014

- Fast, low-latency hardware: Ethernet, SSD
- I/O-intensive applications
- OS on critical (data) path

- Fast, low-latency hardware: Ethernet, SSD
- $\bullet~I/O\text{-intensive}$ applications
- OS on critical (data) path



		Recei	ver Running	CPU Idle		
Network Stack	in out	$\begin{array}{c} 1.26 \\ 1.05 \end{array}$	$(37.6\%)\ (31.3\%)$	$\begin{array}{c} 1.24 \\ 1.42 \end{array}$	$(20.0\%)\(22.9\%)$	
Scheduler		0.17	(5.0%)	2.40	(38.8%)	
Сору	in out	$\begin{array}{c} 0.24 \\ 0.44 \end{array}$	$(7.1\%) \\ (13.2\%)$	$0.25 \\ 0.55$	$(4.0\%)\ (8.9\%)$	
Kernel Crossing	return syscall	$\begin{array}{c} 0.10\\ 0.10\end{array}$	$(2.9\%)\(2.9\%)$	$\begin{array}{c} 0.20 \\ 0.13 \end{array}$	$(3.3\%) \\ (2.1\%)$	
Total		3.36	$(\sigma = 0.66)$	6.19	$(\sigma = 0.82)$	



UDP echo: 1,024 Byte payload, 1,000 samples, times in μs

- $\bullet\,$ Direct application-level access to I/O devices
- Maintain properties of classical server OS: isolation, rate limiting, global naming

- $\bullet\,$ Direct application-level access to I/O devices
- Maintain properties of classical server OS: isolation, rate limiting, global naming
- Use hardware virtualisation mechanisms
 - SR-IOV: virtual PCI devices with dedicated resources
 - Standard for NIC, expected for RAID
 - Memory protection via IOMMU
- Support only policies that hardware can implement efficiently

Arrakis — Hardware Model

Virtual interface cards:

- Managed by control plane
- Dedicated queues and rate limiters
- Virtual network interface card (VNIC)
 - Transmit and receive filters
 - Optional support for offloading
- Virtual storage interface controller (VSIC)
 - Virtual storage areas (VSA)
 - Emulated by dedicated CPU core



- Provides VICs, filters, VSAs, rate specifiers
- IPC endpoints ("doorbells") associated to events on VIC
- Virtual file system interfaces with directories exported by applications



- Applications use (virtual) hardware directly
- Asynchronous operations, doorbells signal completion
- Optimised native interfaces with POSIX-compatibility layer on top
- Network
 - Send/Receive packet to/from queue
 - Extaris: user-level network stack partly based on IwIP
- Storage
 - $\bullet~{\sf Read}/{\sf Write}/{\sf Flush}$ in a VSA; any offset, arbitrary size
 - Caladan: library of persistent data structures (log, queue)

- Implemented by extending Barrelfish
- 6 machine cluster for evaluation
 - 6-core Intel Xeon E5-2430 (Sandy Bridge) @ 2.2 GHz
 - Intel X520 (82599-based) 10Gb Ethernet adapter
 - Intel MegaRAID RS3DC040 RAID controller
 - 100GB Intel DC S3700 SSD
- One machine executes application, others act as clients
- "Tuned" Ubuntu Linux 13.04 (Linux 3.8)

(1,024 Byte payload, 1,000 samples, times in μs)

		Linux				Arrakis				
		Recei	ver Running	CPU Idle		POSIX		Native		
Network stack	in out	$1.26 \\ 1.05$	$(37.6\%) \\ (31.3\%)$	$1.24 \\ 1.42$	(20.0%) (22.9%)	$0.32 \\ 0.27$	(22.3%) (18.7%)	$0.21 \\ 0.17$	(55.3%) (44.7%)	
Scheduler		0.17	(5.0%)	2.40	(38.8%)	_		_		
Сору	in out	$\begin{array}{c} 0.24 \\ 0.44 \end{array}$	(7.1%) (13.2%)	$\begin{array}{c} 0.25 \\ 0.55 \end{array}$	$(4.0\%) \ (8.9\%)$	$0.27 \\ 0.58$	$(18.7\%)\ (40.3\%)$	_		
Kernel crossing	return syscall	$\begin{array}{c} 0.10\\ 0.10\end{array}$	$(2.9\%) \ (2.9\%)$	$\begin{array}{c} 0.20 \\ 0.13 \end{array}$	$(3.3\%) \\ (2.1\%)$	_		_		
Total		3.36	$(\sigma=0.66)$	6.19	$(\sigma=0.82)$	1.44	$(\sigma < 0.01)$	0.38	$(\sigma < 0.01)$	

(1,024 Byte payload, 1,000 samples, times in μs)

		Linux				Arrakis				
		Recei	ver Running	CPU Idle		POSIX		Native		
Network stack	in out	$1.26 \\ 1.05$	$(37.6\%) \\ (31.3\%)$	$1.24 \\ 1.42$	(20.0%) (22.9%)	$0.32 \\ 0.27$	(22.3%) (18.7%)	$0.21 \\ 0.17$	(55.3%) (44.7%)	
Scheduler		0.17	(5.0%)	2.40	(38.8%)	—		—		
Сору	in out	$\begin{array}{c} 0.24 \\ 0.44 \end{array}$	(7.1 %) (13.2 %)	$0.25 \\ 0.55$	$(4.0\%) \\ (8.9\%)$	$\begin{array}{c} 0.27 \\ 0.58 \end{array}$	$(18.7\%)\ (40.3\%)$	_		
Kernel crossing	return syscall	$\begin{array}{c} 0.10\\ 0.10\end{array}$	$(2.9\%)\(2.9\%)$	$\begin{array}{c} 0.20\\ 0.13 \end{array}$	$(3.3\%) \\ (2.1\%)$	_		_		
Total		3.36	$(\sigma=0.66)$	6.19	$(\sigma=0.82)$	1.44	$(\sigma < 0.01)$	0.38	$(\sigma < 0.01)$	



(64k random keys, 1,024 Byte value size, 1,000 samples, times in μs)

		Read	l Hit		Durable Write					
		Linux	А	Arrakis/P		Linux	Arrakis/P			
epoll	2.42	(27.91%)	1.12	(27.52%)	2.64	(1.62%)	1.49	(4.73%)		
recv	0.98	(11.30%)	0.29	(7.13%)	1.55	(0.95%)	0.66	(2.09%)		
Parse Input	0.85	(9.80%)	0.66	(16.22%)	2.34	(1.43%)	1.19	(3.78%)		
Lookup/Set Key	0.10	(1.15%)	0.10	(2.46%)	1.03	(0.63%)	0.43	(1.36%)		
Log Marshaling	_		_		3.64	(2.23%)	2.43	(7.71%)		
write	—		_		6.33	(3.88%)	0.10	(0.32%)		
fsync	—		_		137.84	(84.49%)	24.26	(76.99%)		
Prepare Response	0.60	(6.92%)	0.64	(15.72%)	0.59	(0.36%)	0.10	(0.32%)		
send	3.17	(36.56%)	0.71	(17.44%)	5.06	(3.10%)	0.33	(1.05%)		
Other	0.55	(6.34%)	0.46	(11.30%)	2.12	(1.30%)	0.52	(1.65%)		
Total	8.67	$(\sigma = 2.55)$	4.07	$(\sigma = 0.44)$	163.14	$(\sigma = 13.68)$	31.51	$(\sigma = 1.91)$		

(64k random keys, 1,024 Byte value size, 1,000 samples, times in μs)

		Read	d Hit		Durable Write					
		Linux	Arrakis/P			Linux	Arrakis/P			
epoll	2.42	(27.91%)	1.12	(27.52%)	2.64	(1.62%)	1.49	(4.73%)		
recv	0.98	(11.30%)	0.29	(7.13%)	1.55	(0.95%)	0.66	(2.09%)		
Parse Input	0.85	(9.80%)	0.66	(16.22%)	2.34	(1.43%)	1.19	(3.78%)		
Lookup/Set Key	0.10	(1.15%)	0.10	(2.46%)	1.03	(0.63%)	0.43	(1.36%)		
Log Marshaling	_		_		3.64	(2.23%)	2.43	(7.71%)		
write	—		—		6.33	(3.88%)	0.10	(0.32%)		
fsync	—		—		137.84	(84.49%)	24.26	(76.99%)		
Prepare Response	0.60	(6.92%)	0.64	(15.72%)	0.59	(0.36%)	0.10	(0.32%)		
send	3.17	(36.56%)	0.71	(17.44%)	5.06	(3.10%)	0.33	(1.05%)		
Other	0.55	(6.34%)	0.46	(11.30%)	2.12	(1.30%)	0.52	(1.65%)		
Total	8.67	$(\sigma = 2.55)$	4.07	$(\sigma = 0.44)$	163.14	$(\sigma = 13.68)$	31.51	$(\sigma = 1.91)$		



- Control plane configures hardware to enforce policies
- Direct application access to (virtualised) hardware
- $\bullet\,$ Specialised, application-level I/O stacks
- Significant performance improvements
- Merged back into Barrelfish mainline