The Composite Component-Based Operating System

Gabriel Parmer

Based on research with: Runyu Pan, Yuxin Ren, Phani Kishore Gadepalli, Wenyuan Shao, Qi Wang, Jiguo Song, and many others

The George Washington University gparmer@gwu.edu



Embedded Systems

Predictability Simplicity SWaP-C Multi-tenant Cloud

Performance Isolation Elasticity





Multi-tenant Cloud





Component-based System Design

- Code, data, user-level
- Export APIs (E1 = {fn, ...})
- Explicit dependencies
- Unit of reuse & isolation



- *Minimize* functionality for the *necessary* APIs
- Strong, *fine-grained isolation*

Component-based System Design



Component-based System Design

App3 System composed App1 App2 App3 from components NetStk • Limit *scope* of 12C - compromises NetDrv - faults SHMem Channel - unpredictability MM Sched

Guidance for µ-Kernel Design

"...a concept is tolerated inside the μ -kernel only if moving it outside the kernel, i.e. permitting competing implementations, would **prevent the implementation of the system's required functionality**."

- Liedtke '95

Guidance for µ-Kernel Design

"...a concept is tolerated inside the μ -kernel only if moving it outside the kernel, i.e. permitting competing implementations, would **prevent the implementation of the system's required functionality**."

- Liedtke '95

Composite: push this to the extreme w/ component-defined

- scheduling
- parallel scalability
- concurrency
- capability delegation/revocation

•















Asynchronous IPC



End-to-end timing requires dependency analysis

Synchronous IPC between Threads



End-to-end timing based on kernel mechanims

- In-kernel priority inheritance machinery, ceiling w/ limited prio
- Budget management and experation

Spectrum of IPC Mechanisms











Wang et al. Execution Stack Management for Hard Real-Time Computation in a Component-Based OS, RTSS '11









Parmer et al. Predictable Interrupt Management and Scheduling in the Composite Component-based System, RTSS '08 Gadepalli et al. Temporal Capabilities: Access Control for Time, RTSS '17





Thread Migration schedulers b temporal b, p capabilities

Gadepalli et al. *Temporal Capabilities: Access Control for Time*, RTSS '17 Parmer et al. *HiRes: a System for Predictable Hierarchical Resource Management*, RTAS '11





Gadepalli et al. *Slite: OS Support for Near Zero-Cost, Configurable Scheduling*, RTAS '20







Round-trip IPC	seL4	composite
X86-32 (3.2GHz)	934	741
Cortex a9 (667 GHz, zynq)	630	543

Gadepalli et al. Chaos: a System for Criticality-Aware, Multi-core Coordination, RTAS '19 Samuel Jero et al. Practical Principle of Least Privilege for Secure Embedded Systems, RTAS '21

- Can kernel APIs/implementation
- Limit the scalability of components?
- Cause interference between components?



- Can kernel APIs/implementation
- Limit the scalability of components?
- Cause interference between components?



- Can kernel APIs/implementation
- Limit the scalability of components?
- Cause interference between components?



Component policies for predictability/scalability?

 Kernel must be wait-free w/ component-controlled cacheline contention
 Component-controlled IPI facilities



Gadepalli et al. *Chaos: a System for Criticality-Aware, Multi-core Coordination*, RTAS '19 Wang et al. *Speck: A Kernel for Scalable Predictability*, RTAS '15





Composite Infrastructure

- PoLP-focused RTOS
- NetBSD Rumpkernels
- Xen-like driver domains
- NASA CFS
- Mixed-criticality system orchestration



ucVM Challenges

- MPU-based isolation vs. page-table software abstractions
 - path-compressed radix tries flattened into MPU regions
- Limited # of MPU protected regions
 - Solve memory layout: contiguous memory fits into # regions
 - Treated as dynamic software protection cache
- Efficient coordination and event processing
 - Kernel-bypass + micro-optimization
- ~8 VMs in 128KiB SRAM



ucVM vs. FreeRTOS

Arm Cortex-M7, 216 Mhz

 $A_1 A_2$

FreeRTOS

VMM

Composite uk

A₂

RTOS



Also: Secure Bare-metal interrupts

• Use Trustzone-M as secure kernel bypass for interrupts

Pan et al. Predictable Virtualization on Memory Protection Unitbased Microcontrollers, RTAS '18
Pan et al. MxU: Towards Predictable, Flexible, and Efficient Memory Access Control for the Secure IoT, EMSOFT '19





	Isolation	Scalability	Startup Time	High- performance networking
Processes	*			
Containers	ŗ	ŗ	*	
VMs	\checkmark	*	*	*



	Isolation	Scalability	Startup Time	High- performance networking
Processes	*			
Containers	Ľ	Ľ	*	
VMs	\checkmark	*	*	*
EdgeOS: Feather-weight Processes				



	Isolation	Scalability	Startup Time	High- performance networking
--	-----------	-------------	-----------------	------------------------------------

EdgeOS: Isolation, Predictability, Performance, and Scale

- High speed data movement (10Gbps+) without sacrificing isolation
- Startup > **100X** faster than fork+exec
- Scales to 1000s of services per host = 1 service per user!



EdgeOS: Strong Isolation & Performance



Per-client Service Instantiation

- Docker: the execution time of "docker start"
- Firecracker: the start time of the recommended "hello" image
- Linux: fork() + exec()





Per-client Service Instantiation

- Docker: the execution time of "docker start"
- Firecracker: the start time of the recommended "hello" image
- Linux: fork() + exec()





Per-client Service Instantiation

- Docker: the execution time of "docker start"
- Firecracker: the start time of the recommended "hello" image
- Linux: fork() + exec()





EdgeOS: scalability for TLS proxies





? || /* */

https://github.com/gwsystems/composite https://composite.seas.gwu.edu

