

## Towards Modular Specification and Verification of Concurrent Hypervisor-based Isolation

## Hoang-Hai Dang, David Swasey, Gregory Malecha

BedRock Systems Inc.

## The BedRock Stack

- Microkernel-based systems stack.
- Provides virtualization functionality.
- Kernel security monitors can introspect running guests.
  - Monitoring
  - Protection
  - Remediation

#### **Big Ambitious Goal**

Verify the stack against a strong user-space specification. E.g. "bare metal property".



## The NOVA Microkernel

- Capability-based API supporting dynamic object creation.
- Provides low-level mechanisms for working with system resources, e.g. virtual memory, interrupts, etc.
- Communication using semaphores, IPC, and shared memory.
- Exposes both host threads and "virtual CPUs" (for hardware virtualization)
- User threads are *completely untrusted*.





# Microkernel owns minimum, security-relevant resources

user-owned

µkernel-owned





### NOVA exposes kernel objects and hypercalls and exposes HW-assisted virtualization





#### **Rich Userspace specs and proofs** Outline Part 4. Robust applications. Part 2. Idiomatic specifications Hybrid Clients **Part 3.** Prove *robust safety* from the specification. and application proofs. (verified and untrusted code) System Verification **Robust Safety** (safe against arbitrary clients) (whole-system proofs) Strong specification allows reasoning about the implementation only once. Strong Specification (Separation Logic) Separation logic makes expressing independence easy. **NOVA Implementation** Part 1. NOVA proof against (C++ & ASM) a "precise" specification. **Physical Hardware**



## Fair Semaphores w/ **Timeouts**

State – reflected as separation logic

#### sm.value g n

sm.queue g ls

A queue of blocked ECs

Implement a FIFO discipline to get fairness





#### 4.4.5 Control Semaphore

#### Parameters:

// Semaphore status = ctrl\_sm (SELOB) sm, // Deadline Timeout UINT ticks):

#### Flags:



#### Description:

Prior to the hypercall:

- If D=0 (Up): { PD<sub>CURRENT</sub>, SEL<sub>OB</sub> sm } must refer to a SM Object Capability (CAP<sub>OBJ</sub>) with permission CTRLup.
- If D=1 (Down): { PD<sub>CURRENT</sub>, SEL<sub>OB</sub>, sm } must refer to a SM Object Capability (CAP<sub>OBJec</sub>) with permission CTRL<sub>DN</sub>.

If the hypercall completed successfully:

- If D=0 (Up): if there were ECs blocked on the semaphore, then the microhypervisor has released one of those blocked ECs. Otherwise, the microhypervisor has incremented the semaphore counter. The deadline timeout value and the Z-flag were ignored.
- If D=1 (Down): if the semaphore counter was larger than zero, then the microhypervisor has decremented the semaphore counter (Z=0) or set it to zero (Z=1). Otherwise, the microhypervisor has blocked EC<sub>CURRENT</sub> on the semaphore. If the deadline timeout value was non-zero, EC<sub>CURRENT</sub> unblocks with a timeout status when the architectural timer reaches or exceeds the specified ticks value.

Blocking and releasing of ECs on a semaphore uses the FIFO queueing discipline.

#### SUCCESS

atus:

· The hypercall completed successfully.

#### TIMEOUT

• If D=1: Down operation aborted when the timeout triggered.

#### OVRFLOW

• If D=0: Up operation aborted because the semaphore counter would overflow.

#### BAD\_CAP

• [ PD<sub>CURRENT</sub>, SEL<sub>0BJ</sub> sm ] did not refer to a SM Object Capability (CAP<sub>0BJ</sub>) or that capability had insufficient permissions.

#### BAD CPU

• If D=1 on an interrupt semaphore: Attempt to wait for the interrupt on a different CPU than the CPU to which that interrupt has been routed via assign\_int.

#### 4.4.5 Control Semaphore

#### Parameters:



#### 4.4.5 Control Semaphore



The future is built on

stems Inc

## The NOVA Event Handler

Simplified

```
User & kernel small step semantics
Theorem wp_nova_ec_intro : ∀ ec regs, (inspired by SAIL)
   (\forall \text{ evt regs'}, [| \text{ cpu.step regs evt regs'}] -*
     match evt with
       None => wp_nova_ec regs'
       Some syscall => wp_hypercall ec syscall ...
       Some (mem ..) => wp_mem ..
                                  Address translation and memory access
     end)
                                  (2nd stage page tables managed by NOVA,
                                  reflected in separation logic state)
 ⊢ wp_nova_ec ec regs.
```



## **Starting NOVA (Simplified)**

Parametric over *any* startup program.





Input

## Starting NOVA (Simplified) Userspace









## System Refinement

Establish properties for the applications that run on top of NOVA.

- Use CaReSL-style techniques to prove refinement using ghost state and invariants.
- Extract the end-to-end proof (independent of SL) using Iris adequacy.





## System Refinement

Establish properties for the applications that run on top of NOVA.

- Use CaReSL-style techniques to prove refinement using ghost state and invariants.
- Extract the end-to-end proof Specification. Specification.
- spec init ⊢ nova\_state -\* elf user.bin -\*
  wp\_nova\_ec ...
  spec init ⊢ nova\_ok Use the proof of NOVA.
  spec init ⊢ physical\_state -\* wp\_arm boot
  Use adequacy to extract a standard
  operational refinement proof.



## System Refinement

Establish properties for the applications that run on top of NOVA.

- Use CaReSL-style techniques to prove refinement using ghost state and invariants.
- Extract the end-to-end proof (independent of SL) using Iris adequacy.









## Proving Robust Safety

Want to show that NOVA is safe for **any** userspace code.

▼ startup_image, Veed to show that this assumption is trivial.
<pre>(∀ root_pd root_ec root_sc, pd.mem root_pd startup_image -* initial_pd state root_pd root_ec root_sc -* memory -** wp_nova_ec root_ec (startup_regs))</pre>
<pre>     NOVA_loaded -*     elf startup_image -*     memory -*    *     wp_arm_el2 boot_regs</pre>



## Proving Robust Safety from the Spec

Express the high-low state distinction within separation logic.

- Invariants allow flexible, concurrent sharing.
- Existential quantifiers express low-integrity

 ∃ objs, let valid o := o \in objs in rs-inv [\*list] o \in objs,
 ☐ v es, sm.value γ v \*\* sm.queue γ ecs \*\* [| Forall valid es |]

"Low integrity" invariant for semaphores. Ownership exists, values are minimally constrained.

Properties for other kernel object types.

\*\* all schedulable ECs are **valid** \*\* all interrupts bound to **valid** interrupt SMs

Need to Show

Userspace Resources |-- |={T}=> inv rs-inv inv rs-inv |-- wp\_nova\_ec boot\_ec boot\_regs

Over arbitrary user binaries.

Systems Inc





Userspace Robust Safety

State may change integrity levels throughout the execution of the program.

- High integrity state shared with untrusted code.
- Low integrity state revoked from untrusted code.

Revocation is generally difficult and requires tight reasoning about confinement and visibility.

Precise specifications support endorsement (low -> high) because they decouple state from policy.

ROCK

Provide mechanisms (assertions), not policies (invariants, quantifiers). Support a "data life cycle".

- Entire data lifecycle is provable using the strong specification.
- 1. Create state. (high)
- 2. Configure state. (high)
- 3. Share permissions with untrusted code? (low)
- 4. Revoke state. (high)
- 5. Destroy state. (high)

Can share limited permissions, e.g. only up a semaphore, only call a portal, only read a page, etc.



#### **Rich Userspace specs and proofs** Outline Part 4. Robust applications. Part 2. Idiomatic specifications Hybrid Clients **Part 3.** Prove *robust safety* from the specification. and application proofs. (verified and untrusted code) System Verification **Robust Safety** (safe against arbitrary clients) (whole-system proofs) Strong specification allows reasoning about the implementation only once. Strong Specification (Separation Logic) Separation logic makes expressing independence easy. **NOVA Implementation** Part 1. NOVA proof against (C++ & ASM) a "precise" specification. **Physical Hardware**

