

## Modular, Full-Stack Verification

**Gregory Malecha**, Hoang-Hai Dang, Paolo G. Giarrusso, Simon Hudon, Jan-Oliver Kaiser (BlueRock Security) David Swasey (Riverside Research)

gregory@bluerock.io

We can *specify real, concurrent systems* in ways that *naturally align with systems design principles*.



Full verification of real systems code is tractable.

## Modular, Full-Stack Verification

**Gregory Malecha**, Hoang-Hai Dang, Paolo G. Giarrusso, Simon Hudon, Jan-Oliver Kaiser (BlueRock Security) David Swasey (Riverside Research)

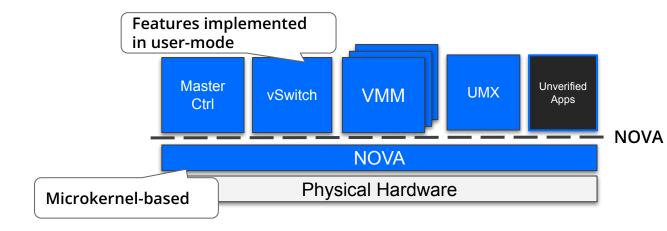
gregory@bluerock.io

## The Challenge A Virtualization Platform



## The Challenge A Virtualization Platform

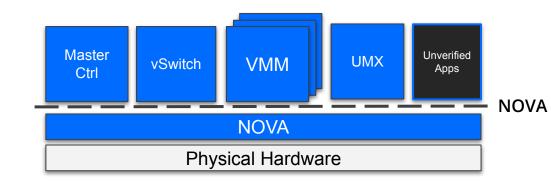




## The Challenge A Virtualization Platform

- 1. **Real System** concurrency, dynamic resource sharing, real systems engineers (+mainstream PL, e.g. C++).
- 2. **Two-sided Specification** single specification used by both the implementation and the clients.
- 3. **Modular Verification** aligned with the architecture of the actual systems engineers.
- 4. **Robust** support inter-operation with unverified code

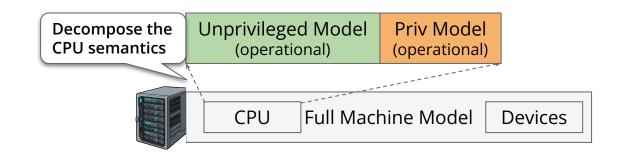




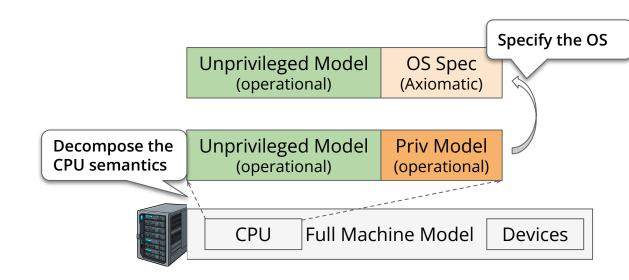




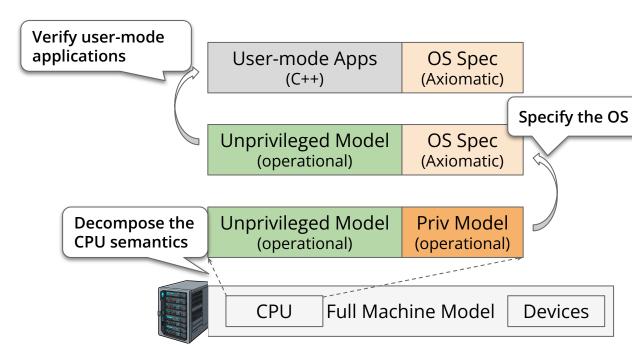














User-mode Apps (C++) OS Spec (Axiomatic)

Unprivileged Model (operational)

OS Spec (Axiomatic)

Unprivileged Model (operational)

Priv Model (operational)

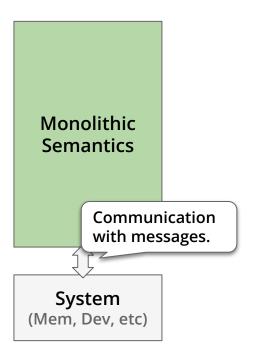


CPU

Full Machine Model

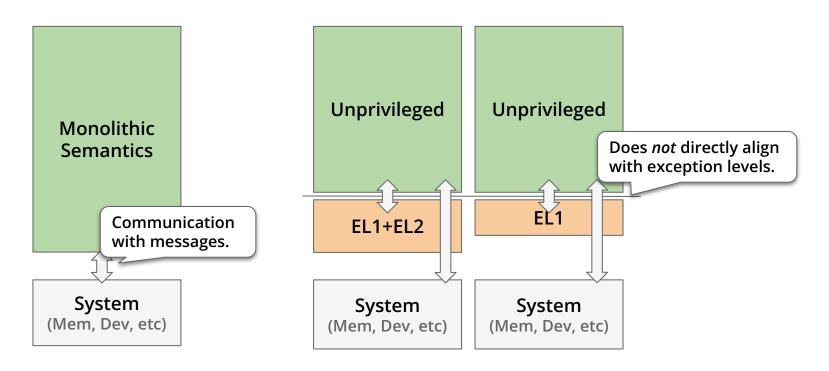
Devices

## **Decomposing CPU Semantics**

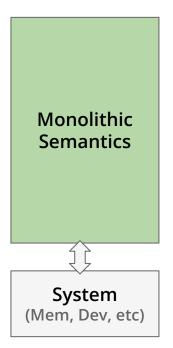


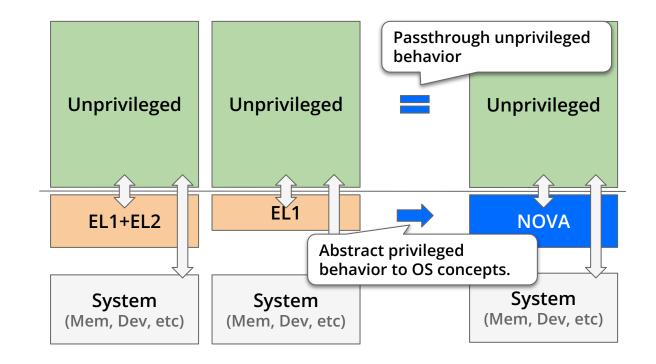


## **Decomposing CPU Semantics**



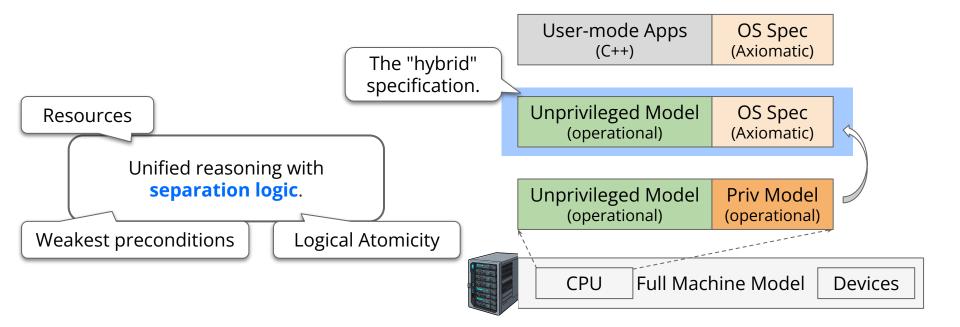
## **Decomposing CPU Semantics**





# The Full-Stack Proof Specify the OS



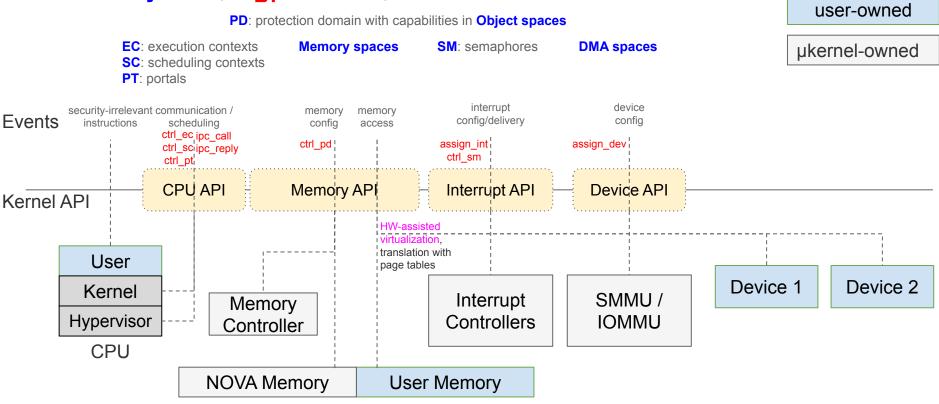


### The NOVA Machine

```
This is a "thread"-local
                          reasoning principle!
Theorem wp_nova_ec_intro : ∀ ec regs,
   (∀ evt regs',
                                            _ *
        cpu.step regs evt regs'
                                                                       Unprivileged
     match evt with
                                           Hypercall behavior
       None => wp_nova_ec regs'
       Some syscall => wp_hypercall ec ..
       Some (mem vaddr) => wp_mem vaddr
                                                                          NOVA
                                 Address translation specified
     end)
                                 using NOVA state.
   wp_nova_ec ec regs.
                                                                         System
                                               <del>Jy3teiii</del>
(Mem, Dev, etc)
                                                                       (Mem, Dev, etc)
                            (Mem, Dev, etc)
                                            (Mem, Dev, etc)
```

## **NOVA's Capability API:**

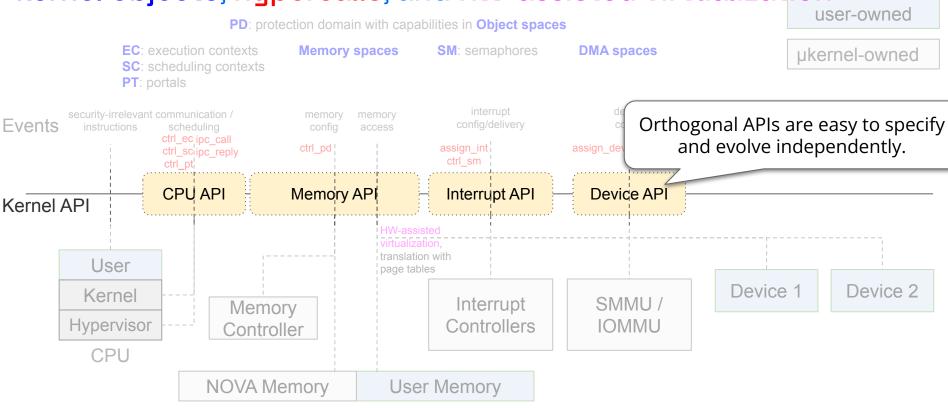
kernel objects, hypercalls, and HW-assisted virtualization





## **NOVA's Capability API:**

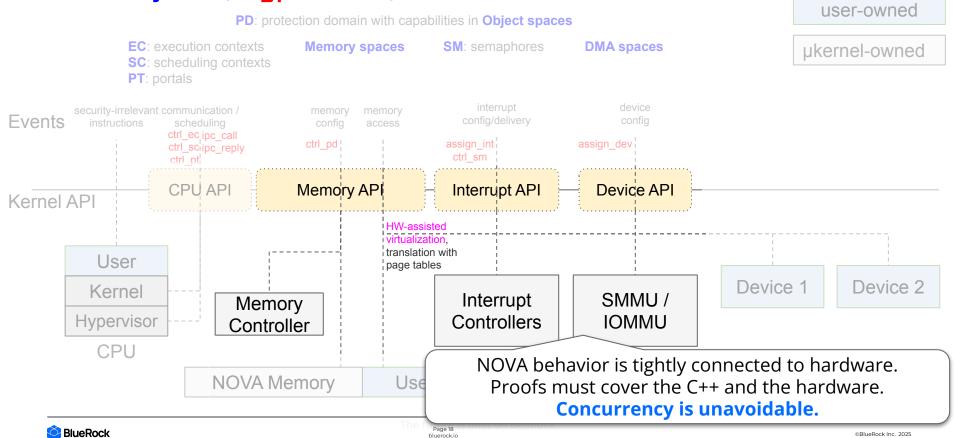
kernel objects, hypercalls, and HW-assisted virtualization





## **NOVA's Capability API:**

kernel objects, hypercalls, and HW-assisted virtualization



## Fair Semaphores

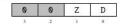


#### 4.4.5 Control Semaphore

#### Parameters:

```
status = ctrl_sm (SELOBJ sm, // Semaphore
UINT ticks); // Deadline Timeout
```

#### Flags:



#### Description:

Prior to the hypercall:

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

#### Status:

#### SUCCESS

· 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

 PDCURRENT, SELOBJ Sm } did not refer to a SM Object Capability (CAPOBJOR) 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.

### Fair Semaphores State

#### Parameters:

```
status = ctrl_sm (SELOBJ sm, // Semaphore
UINT ticks); // Deadline Timeout
```

## A minimal piece of independent state and an ownership discipline/"protocol"

### Reflected as separation logic resources

Selectors map to capabilities

```
pd;sel → (obj, perm)
```

The counter of a semaphore

```
sm.counter obj n
```

• The queue of blocked threads

```
sm.queue obj Is
```

#### Prior to the hypercall:

4.4.5 Control Semaphore

- If D=0 (Up): { PD<sub>CURRENT</sub>, SEL<sub>OBJ</sub> sm } must refer to a SM Object Capability (CAP<sub>OBJ<sub>MR</sub></sub>) with permission CTR<sub>Lup</sub>.
- If D=1 (Down): { PD<sub>CURRENT</sub>, SEL<sub>OBJ</sub> sm } must refer to a SM Object Capability (CAP<sub>OBJ<sub>SS</sub></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.

#### Status:

#### SUCCESS

· 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

 PDCURFENT, SELOBJ sm } did not refer to a SM Object Capability (CAPOBJox) 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.



## Fair Semaphores Behavior

### Reflected as separation logic resources

Selectors map to capabilities

The counter of a semaphore

```
sm.counter obj n
```

• The queue of blocked threads

sm.queue obj ls

#### 4.4.5 Control Semaphore

#### Parameters:

```
status = ctrl_sm (SELOBJ sm, // Semaphore
UINT ticks); // Deadline Timeout
```

#### Flags:



#### Description:

Prior to the hypercall:

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

a If D=1 (Down): if the semanhore counter was larger than zero, then the microhypervisor has

```
let (Some obj) :=
  resolve_sm ec.pd {UP} sel else K BAD_CAP;
sm.incr obj K
```

⊢ wp\_hypercall ec ctrl\_sm(0, sel, \_) K

If D=1: Down operation aborted when the time

#### Continuation.

#### OVRFLOW

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

#### BAD CAP

 PDCURRENT, SELOBJ Sm } did not refer to a SM Object Capability (CAPOBJON) 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.



Sequencing

#### 4.4.5 Control Semaphore

status = ctrl\_sm (SELOB) sm,

#### Parameters:

```
UINT ticks):
                                                                      // Deadline Timeout
Fair Semaphores
                                             Exchange this...
Behavi
         sm.incr g K :=
            AU<< 1 n sm.counter g
                                                                ...for this
Reflected
                                                                                    (CAPOBJON) with
                << sm.counter g (cap</pre>
                                                                                    (CAPOBJew) with
    Sele
                             (if overflows (n+1)
                                                                                    as released one
                                                                                    e counter. The
     po
                                then OVERFLOW else SUCCESS)>>
         Afterwards, continue
                                          >solve_sm ec.pd {UP} sel else K BAD_CAP;
      sm.counter obj n
                                       sm.incr obj K
                                   ⊢ wp_hypercall ec ctrl_sm(0, sel, _) K
    The queue of blocked threads
```

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\_CIRRENT, SEL\_08] sm } did not refer to a SM Object Capability (CAP\_08Jsr) or that capability had insufficient permissions.

// Semaphore

#### 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.



sm.queue obj ls

## Fair Semaphores Behavior

#### 4.4.5 Control Semaphore

#### Parameters:

```
status = ctrl_sm (SEL<sub>OBJ</sub> sm, // Semaphore
UINT ticks); // Deadline Timeout
```

#### Flags:

```
0 0 Z D
```

Description:

Reflected as **separation logic resource** resolve\_sm pd sel K :=

Selectors map to capabilities

```
pd;sel → (obj, perm)
```

The counter of Two atomic steps.

sm.counter obj n

• The queue of blocked threads

```
sm.queue obj Is
```

```
AU<<∃ obj,perm. pd;sel→(obj,perm)>> @NOVA

<<ec.pd;sel→(obj,perm),

COMM K (if {UP} ⊆ perm

then Some obj else None)>>
```

```
let (Some obj) :=
  resolve_sm ec.pd {UP} sel else K BAD_CAP;
sm.incr obj K
```

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.



## Fair Semaphores Behavior

### Reflected as separation logic resource

Selectors map to capa

pd;sel → (obj, pe

• The counter of

sm.counter

The queue of blocked

sm.queue obj Is



Parameters:

Flags:

4.4.5 Control Semaphore

status = ctrl\_sm (SELOB) sm,

UINT ticks):

>> @NOVA

nen Some obj

THE CALL PAD CAP;

// Semaphore

// Deadline Timeout

- 1. Connect the abstract state to C++ state
- 2. Apply the linearization points when updating the state.

```
AU<<∃ n. sm.counter g n>> @NOVA
<<sm.counter g (cap (n+1)),
COMM K (if overflows (n+1)
then OVERFLOW else SUCCESS)>>
```

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.



## Spawning Threads Concurrent Behavior

Thread creation occurs when scheduling contexts are bound to global execution contexts.

The proof obligation for a thread creation is expressed as a WP for the new thread.

#### 5.3.3 Create Scheduling Context

#### Parameters:

#### Flags:



#### Description:

Creates a new Scheduling Context (SC).

Prior to the hypercall:

- SPC<sub>OBJCURRENT</sub>[sel] must refer to a Null Capability (CAP<sub>0</sub>).
- SPC<sub>OBJ<sub>CURRENT</sub></sub>[pd] must refer to a PD Capability (CAP<sub>OBJ<sub>PD</sub></sub>) with permission SC.
- SPC<sub>OBJ<sub>CURRENT</sub></sub>[ec] must refer to an EC Capability (CAP<sub>OBJ<sub>EC</sub></sub>) with permission BIND<sub>SC</sub>.

If the hypercall completed successfully:

- · A new Scheduling Context (SC) has been created.
- The created SC is bound to the EC referred to by SPC<sub>OBJ<sub>CURRENT</sub></sub>[ec] on the CPU of that EC, with its scheduling parameters set according to scd.
- The resources for the created SC were accounted to the PD referred to by SPCOBLINGER [pd].

```
let g_pd := resolve_pd ec.pd sel_pd {SC} Q in
let g_ec := resolve_ec ec.pd sel_ec {EC_BIND_SC} Q in
let g_sc := sc.create ec.pd sel_sc Q in
let run := ec.bind g_ec g_sc in

(if run then wp_nova_ec g_ec else emp)* K SUCCESS
```

⊢ wp\_hypercall ec create\_sc(sel\_sc, sel\_pd, sel\_ec, scd) K

#### BAD PAR

At least one SCD field in scd was invalid.

MEM\_OBJ



### Spawning Threads Concurrent Behavior

Thread creation occurs when scheduling contexts are bound to global execution contexts.

The proof obligation for a thread creation is expressed as a WP for the new thread.

#### 5.3.3 Create Scheduling Context

#### Parameters:

#### Flags:



#### Description:

Creates a new Scheduling Context (SC).

Prior to the hypercall:

- SPC<sub>OBJCURRENT</sub> [sel] must refer to a Null Capability (CAP<sub>0</sub>).
- SPC<sub>OBJCURRENT</sub>[pd] must refer to a PD Capability (CAP<sub>OBJPD</sub>) with permission SC.
- SPC<sub>OBJCURRENT</sub> [ec] must refer to an EC Capability (CAP<sub>OBJEC</sub>) with permission BIND<sub>SC</sub>.

If the hypercall completed successfully:

- A new Scheduling Context (SC) has been created.
- The created SC is bound to the EC referred to by SPCOBICTURED [ec] on the CPU of that EC, with its scheduling parameters set according to scd.
- The resources for the created SC were accounted to the PD referred to by SPCOBLORDOR [pd]

let g\_pd := resolve\_pd ec.pd sel\_pd {SC} Q in -- rocalve\_ec ec.pd sel\_ec {EC\_B\_IND Prove an extra WP for Continue running the reate ec.pd sel\_sc Q in current thread. the new thread. hd g\_ec g\_sc in

(if run then wp\_nova\_ec g\_ec else emp)\* K SUCCESS

⊢ wp\_hypercall ec create\_sc(sel\_sc, sel\_pd, sel\_ec, scd) K

#### BAD PAR

· At least one SCD field in scd was invalid.

#### MEM OBJ





## Spawning Threads Concurrent Behavior

Thread creation occurs when scheduling contexts are bound to global execution contexts.

The proof obligation for a thread creation is expressed as a WP for the new thread.

#### 5.3.3 Create Scheduling Context

#### Parameters:

#### Flags:



#### Description:

Creates a new Scheduling Context (SC).

Prior to the hypercall:

- SPC<sub>OBJCURRENT</sub>[sel] must refer to a Null Capability (CAP<sub>0</sub>).
- SPC<sub>OBJ<sub>CURRENT</sub></sub>[pd] must refer to a PD Capability (CAP<sub>OBJ<sub>PD</sub></sub>) with permission SC.
- SPC<sub>OBJ<sub>CURRENT</sub></sub>[ec] must refer to an EC Capability (CAP<sub>OBJ<sub>EC</sub></sub>) with permission BIND<sub>SC</sub>.

If the hypercall completed successfully:

- · A new Scheduling Context (SC) has been created.
- . The created SC is bound to the EC referred to by SPCOBJCURGENT [ec] on the CPU of that EC, with its

Need logical atomicity to make this complete.

PD referred to by SPC<sub>OBJCURRENT</sub>[pd].

CC)

in

Independence is the core of concurrency!

Continue running the current thread.

(if run then <a href="wp\_nova\_ec g\_ec">wp\_ec</a> else emp) \* K <a href="SUCCESS">SUCCESS</a>

⊢ wp\_hypercall ec create\_sc(sel\_sc, sel\_pd, sel\_ec, scd) K

#### BAD\_PAR

At least one SCD field in scd was invalid.

#### MEM\_OBJ





let g\_pd := resolve\_pd

Prove an extra WP for

the new thread.

# The Full-Stack Proof User-mode Verification

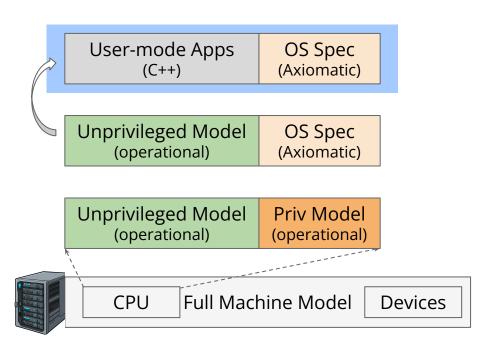
- Device drivers.
  - PL011
  - o (Simple-mode) Zynqmp
- Terminal multiplexor
- VIRTIO-based layer 2 switch.
- Virtual Machine Monitor

**Unified reasoning** with separation logic.

Highly reusable proofs

Malecha, et. al. **Developing with Formal Methods at BedRock Systems**. IEEE S&P Giarrusso, et. al. **Verifying a Virtual Machine Monitor**. BlueRock Technical Report







# Uniform Reasoning in Separation Logic

```
Definition portal_spec_cpu
                                   User-mode entry point (from
    (exit reason : evt.t nova.)
    (yg : GlobalRoundupInfo.ghos
                                      hardware virtualization)
    (cy : bm.cpu.Name → Vcpu.gna
    (shareds : bm.cpu.Name → Vcpu.shared t)
    (bm cid : bm.cpu.Name)
    : WpSpec_cpp :=
  \let yvcpu := cy bm_cid
  \pre{msea : Vcpu.seq_t} vcpu |-> Vcpu.seqR (cQp.m 1) bm_cid yvcpu msea
 \let pd := (Model.Cpu.pd_name (Vcpu.base_basev vvcpu))
  \prepost{yboard qboard mboard} Model.Board.P pd qboard yboard mboard
  \require yboard .^ Model.BoardImpl.Model._stage = Model.Board.SteadyState
  \prepost{q_vcpus} _global "Model::vcpus" |-> Vcpu.vcpusR cy shareds q_vcpus yg
  \prepost{vcpu_regs : arch.regs_t guest}
   Model.Cpu.ecRegs (Vcpu.base_basey yvcpu) vcpu_regs
  \pre{bm_state : lts_state (bm.cpu.lts _)} bm.cpu.core bm_cid bm_state
  \require{decode_assist} cpu.pendingInstruction bm_state decode_assist
  \require{fault_regs} VCPU.FaultRegs.reg_accessor_encodes_decodeassist fault_regs
                         exit_reason decode_assist
  \require CPU.related_full_states bm_state vcpu_regs
  (* ... other ownership ... *)
```

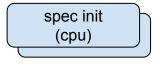
# Uniform Reasoning in Separation Logic

BlueRock

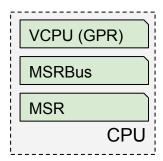
```
User-mode entry point (from
                        (exit reason : evt.t nova.)
                        (yg : GlobalRoundupInfo.ghos
                                                          hardware virtualization)
                        (cy : bm.cpu.Name → Vcpu.gna
                        (shareds : bm.cpu.Name → Vcpu.shared t)
                        (bm cid : bm.cpu.Name)
                        : WpSpec_cpp :=
                      (** Sequential ownership of current VCPU -- C++/NOVA *)
                      \let yvcpu := cy bm_cid
                       \pre{msea : Vcpu.seq_t} vcpu |-> Vcpu.seqR (cQp.m 1) bm_cid yvcpu msea
        C++ state.
                         * Shared handle to the <<Board>> -- C++/NOVA *)
                       (let pd := (Model.Cpu.pd_name (Vcpu.base_basey yvcpu))
  and some
                       \prepost{yboard qboard mboard} Model.Board.P pd qboard yboard mboard
 NOVA state
                      \require yboard .^ Model.BoardImpl.Model._stage = Model.Board.SteadyState
                      (** Shared handles to other VCPUs -- C++/NOVA *)
                      \prepost{q_vcpus} _global "Model::vcpus" |-> Vcpu.vcpusR cy shareds q_vcpus yg
    NOVA State
                      (** VCPU registers -- NOVA *)
                      prepost{vcpu_regs : arch.regs_t guest}
                        Model.Cpu.ecRegs (Vcpu.base_basey yvcpu) vcpu_regs
                       (** Specification state of the Guest -- spec *)
Specification state
                       pre{bm_state : lts_state (bm.cpu.lts _)} bm.cpu.core bm_cid bm_state
                       require{decode_assist} cpu.pendingInstruction bm_state decode_assist
                       \require{fault_regs} VCPU.FaultRegs.reg_accessor_encodes_decodeassist fault_regs
                                              exit reason decode assist
                      (** Spec & Implementation state are related! *)
                       require CPU.related_full_states bm_state vcpu_regs
       Simulation
                         ... other ownership ... *)
```

Definition portal\_spec\_cpu

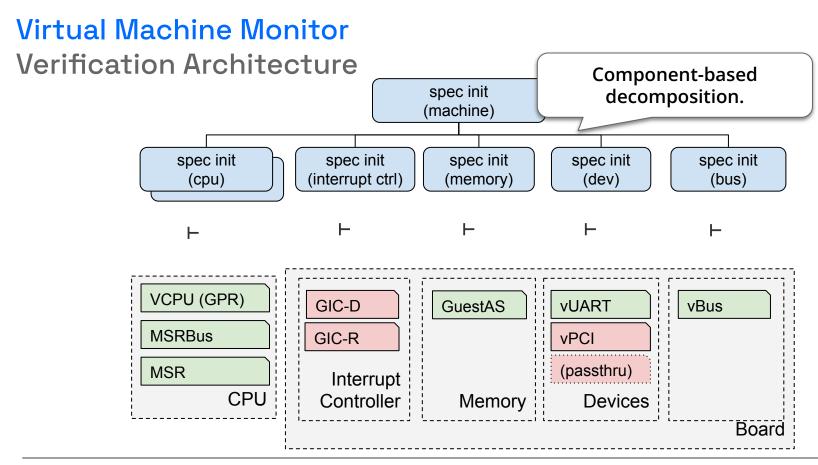
## Virtual Machine Monitor Verification Architecture



 $\vdash$ 







## **Takeaways**

- Separation logic is a powerful formalism for reasoning about dynamic systems.
- It can be used to specify complex, concurrent behavior.
- It can be used to modularly verify complex, concurrent behavior.
- It works for the (good) *code that you have right now*.

#### The Future

- Specify your next system using SL!
- Verify the implementation!

#### **Future Research**

- Refine low-level hardware models
  - Security reasoning.
  - Weak memory.
- Liveness verification.
  - o Availability.
- Hyper-properties.
  - o Confidentiality.
  - Isolation.