

Mechanized Verification of Preemptive OS Kernels

Xinyu Feng

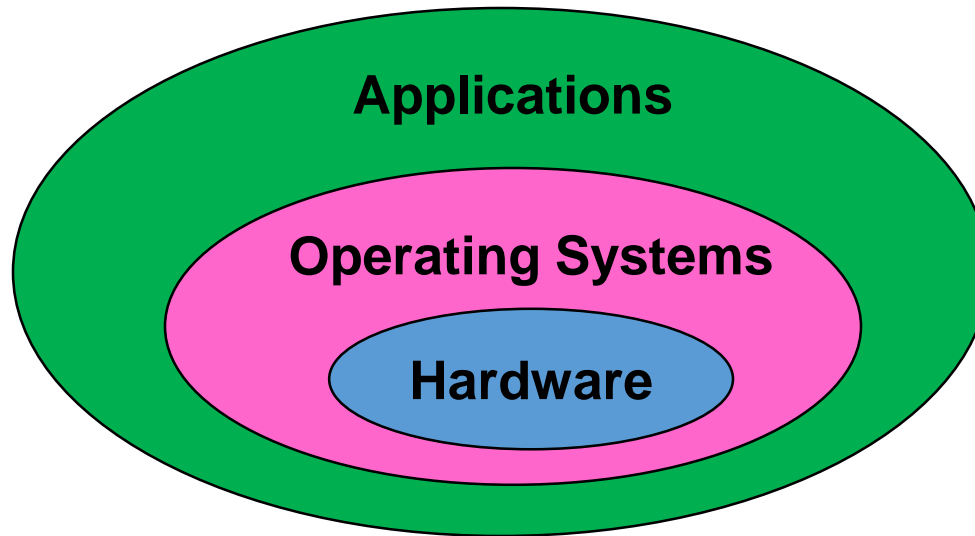
University of Science and Technology of China

Why OS Kernel Verification?



Computer Systems

Why OS Kernel Verification?



**Correctness of OS is crucial for safety
and security of the whole system**

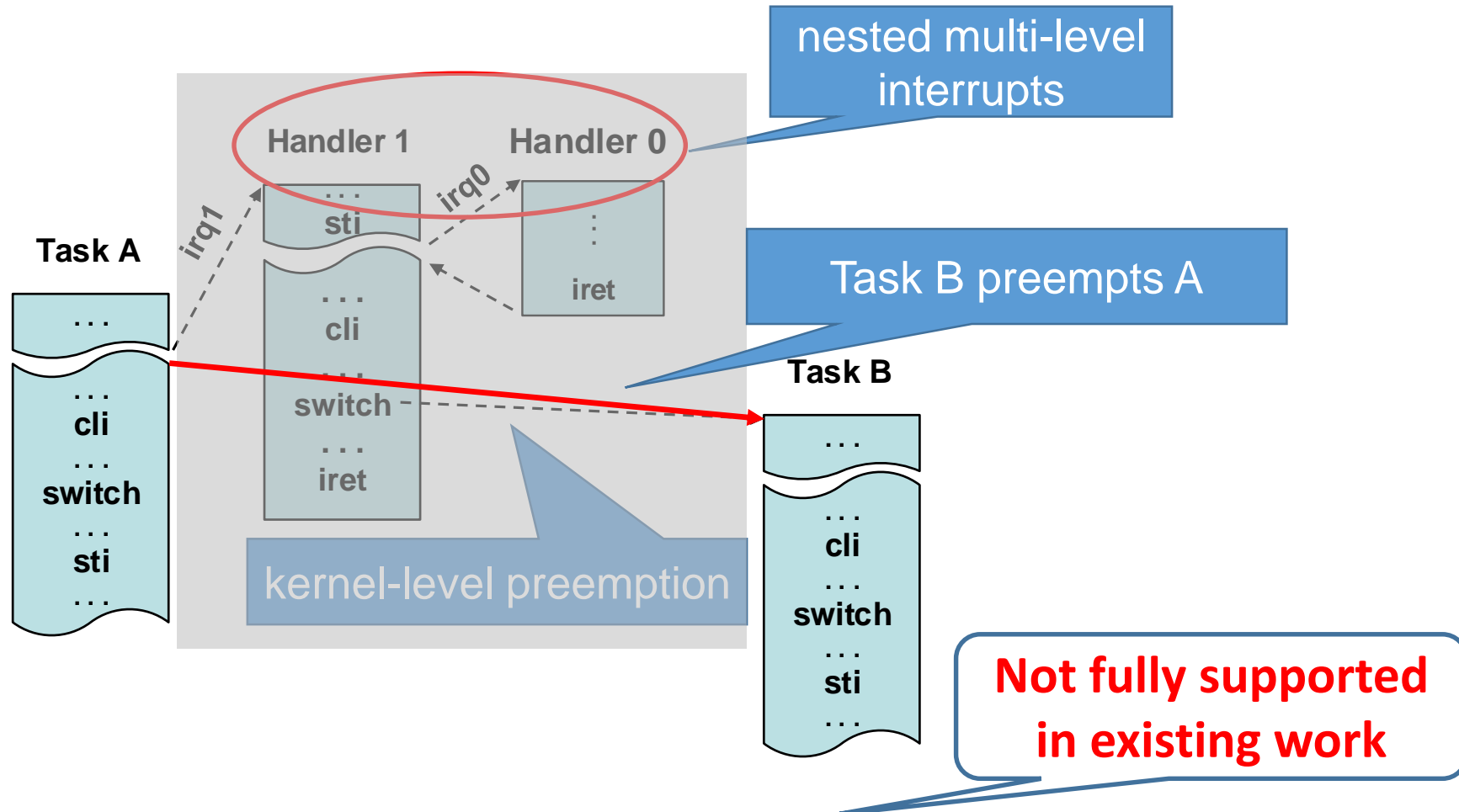
Why OS Kernel Verification?

- Fundamental, but also **simpler to verify!**
(comparing to applications)
 - Less domain knowledge required
 - every programmer knows OS
 - Stable specifications
 - Slow evolution
 - Specs validated by application-level verification

OS Kernel Verification: Challenges

- Low-level programs
 - C + inline assembly, interrupts, task management, ...
- Larger code base (than algorithm verification)
- Code at different abstraction layers
 - E.g., threads vs. schedulers
- Involves both libraries (sys. calls) and runtime (scheduler)
 - What is a proper specification?
- Rich concurrency
 - Multi-tasking, multi-core, interrupts

Preemption and nested interrupts



Preemptions and multi-level interrupts are crucial for real-time systems.

They also make system highly concurrent and complex

Concurrency & Preemption in Previous work

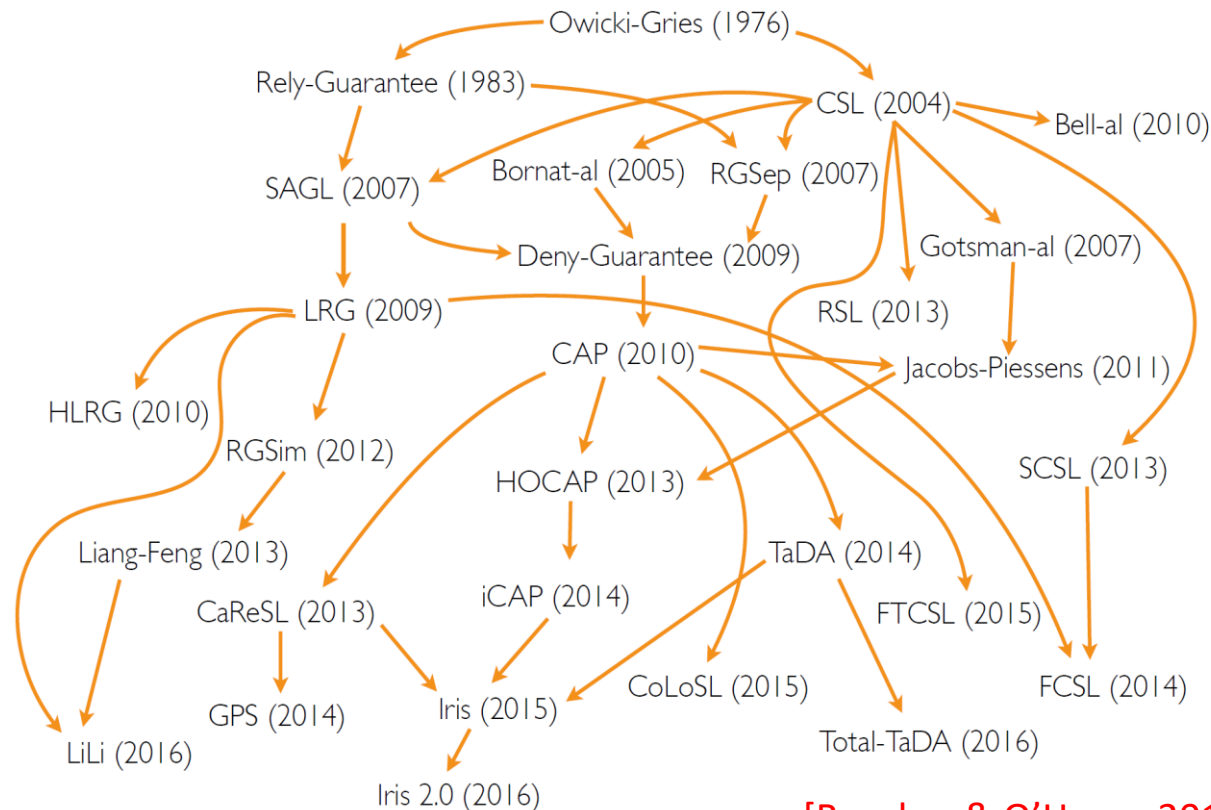
- seL4 [Klein et al. 2009 ...]
 - Mostly sequential
 - Limited support of interrupts at fixed program points
- Verisoft [Rieden et al. 2007 ...]
 - Kernel is sequential
- Verve [Yang & Hawblitzel. PLDI 2010]
 - Allows preemption, but no nested interrupts
 - Mostly about safety, limited functionality verification
- CertiKOS [Gu et al. 2015, Chen et al. 2016, Gu et al. 2016]
 - Evolving: sequential → limited interrupts → multicore
 - Still no preemption

Concurrency & Preemption in Previous work (2)

- eChronos OS [\[Andronick et al. 2015, 2016\]](#)
 - Supports preemption and nested interrupts
 - But verification at the model level only
 - Verifies scheduling invariants, no API correctness

Challenges for Verifying **Preemptive** OS Kernels

- Verifying concurrent programs is difficult
 - Non-deterministic interleaving



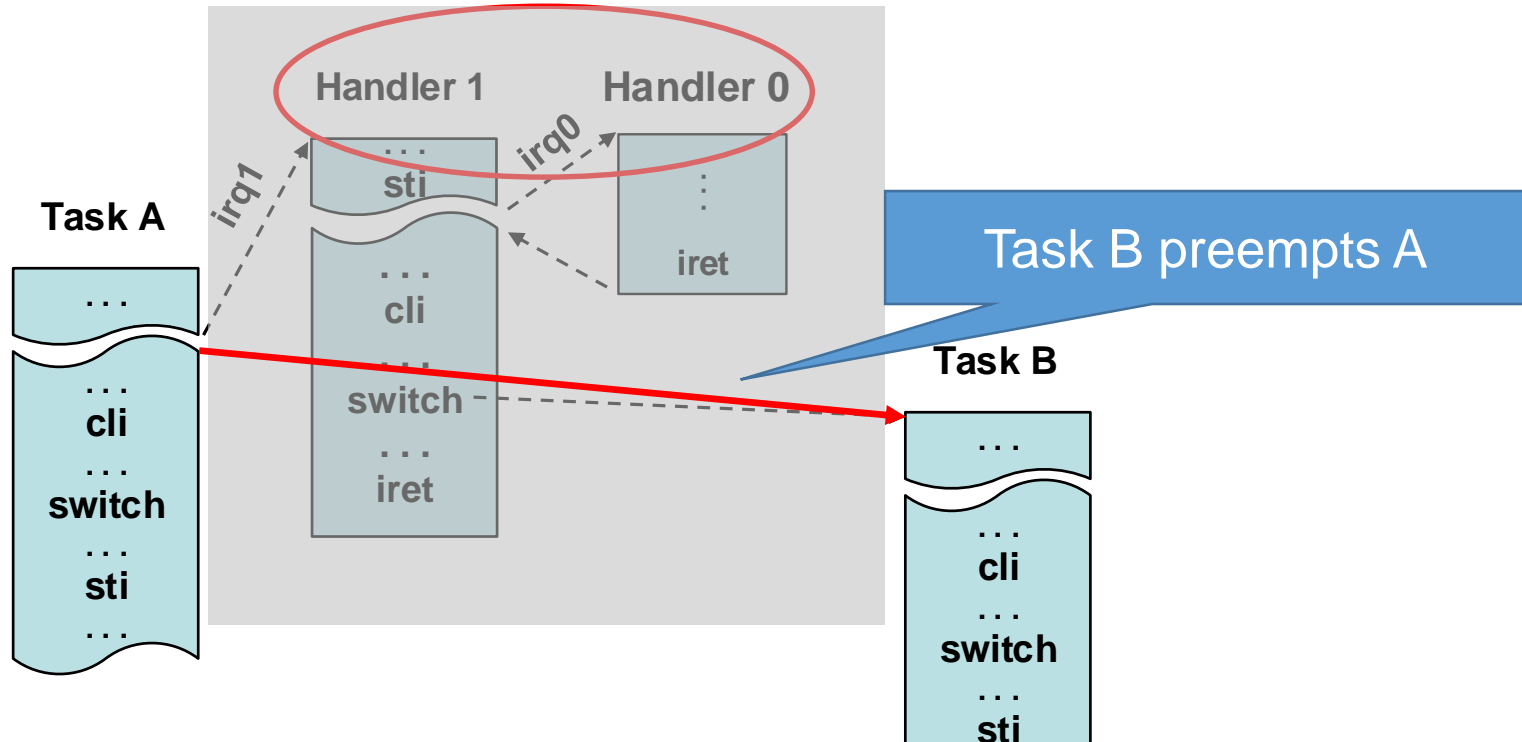
[Brookes & O'Hearn 2016], courtesy of Ilya Sergey

Challenges for Verifying Preemptive OS Kernels

- Verifying concurrent programs is difficult
- Verifying concurrent kernels is even more challenging
 - More difficult to establish **refinement** with concurrency
 - Theories not fully developed until recently
[Turon et al. POPL'13, ICFP'13] [Liang et al. PLDI'13, CSL-LICS'14]
 - Kernel-level preemption can be more complex than multi-tasking/multi-processor concurrency

A natural correctness spec. for OS kernels

Kernel-level preemption can be more complex than multi-tasking/multi-processor concurrency



Interrupt management is now a verification target:
lower abstraction layer and non-uniform concurrency model

More low-level details:

e.g., can context switch only when there are no nested interrupts

This talk

- **Verification framework for preemptive OS kernels**
 - Refinement reasoning about concurrent kernels
 - Multi-Level nested interrupts and preemption
- **Verification of key modules of a commercial OS kernel $\mu\text{C}/\text{OS-II}$ in Coq**



The **first** mechanized verification of
a commercial **preemptive** OS kernel.



[Xu et al. CAV'16]

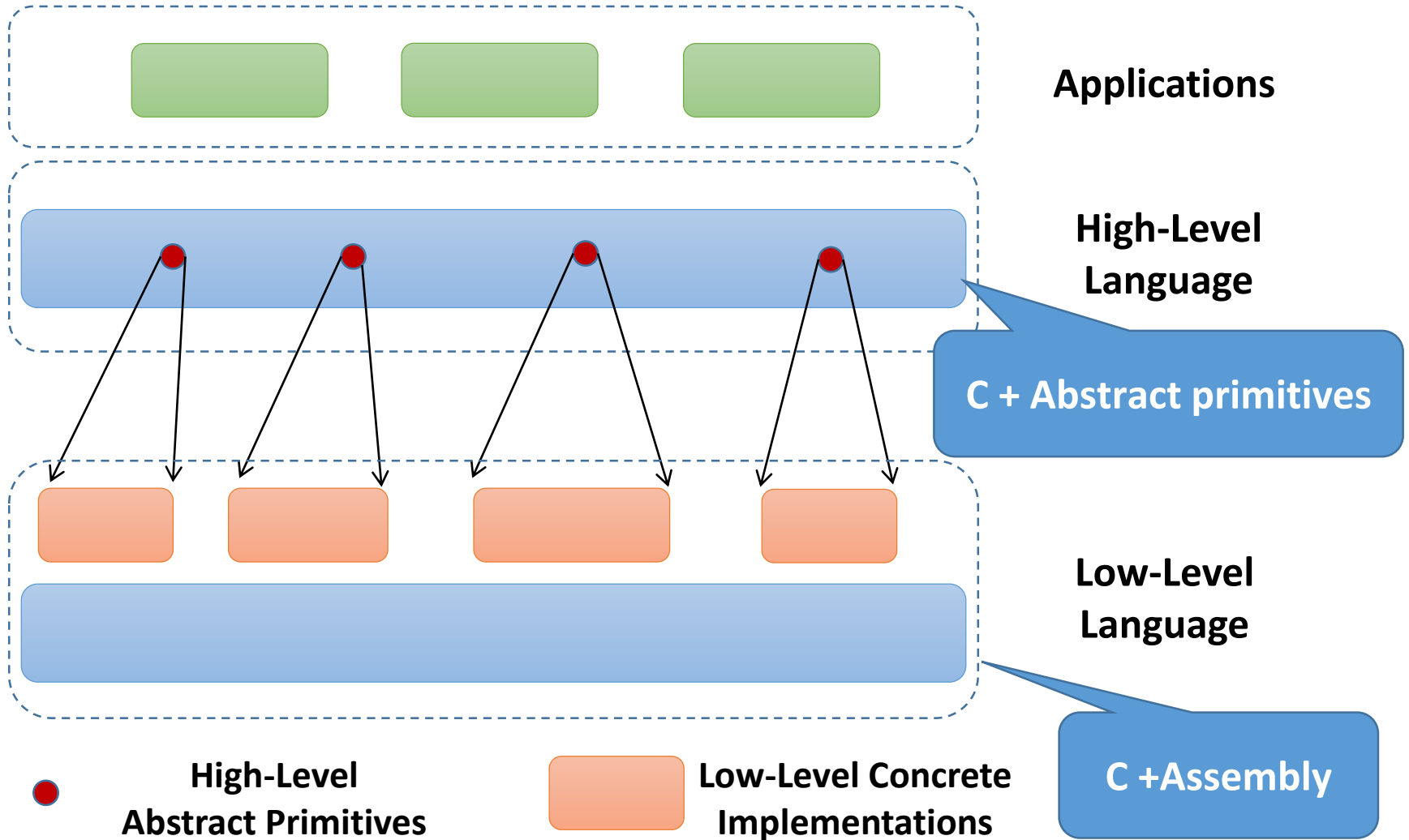
Outline

- OS Correctness Specification
- Verification Framework
 - System modeling
 - CSL-R: Program logic for refinement & multi-level interrupts
 - Coq tactics
- Verifying $\mu\text{C}/\text{OS-II}$

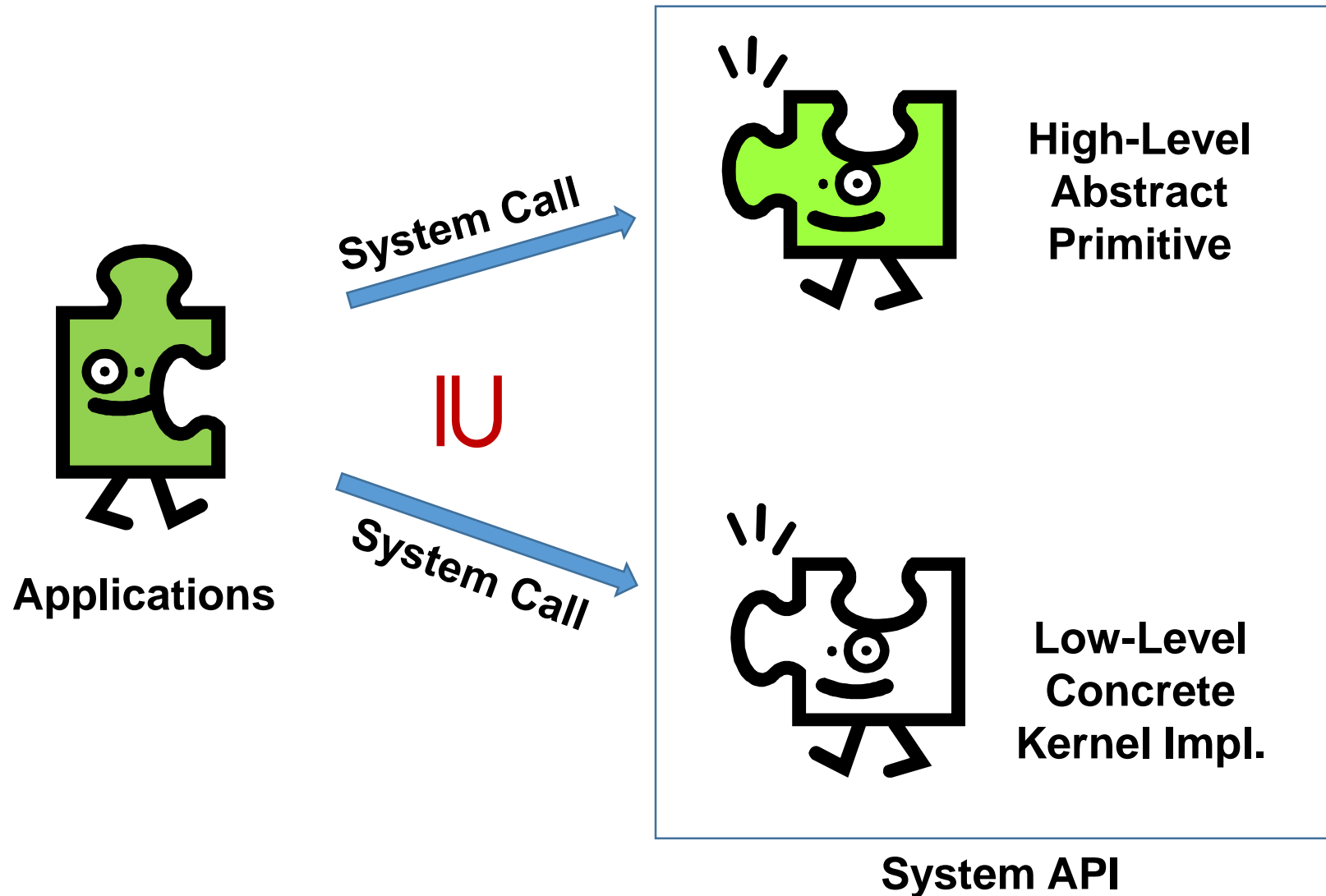
OS Correctness

- **OS provides abstraction for programmers**
 - Hides details of the underlying hardware
 - Provides an abstract programming model
- **OS Correctness : refinement** between high-level abstraction and low-level concrete implementation

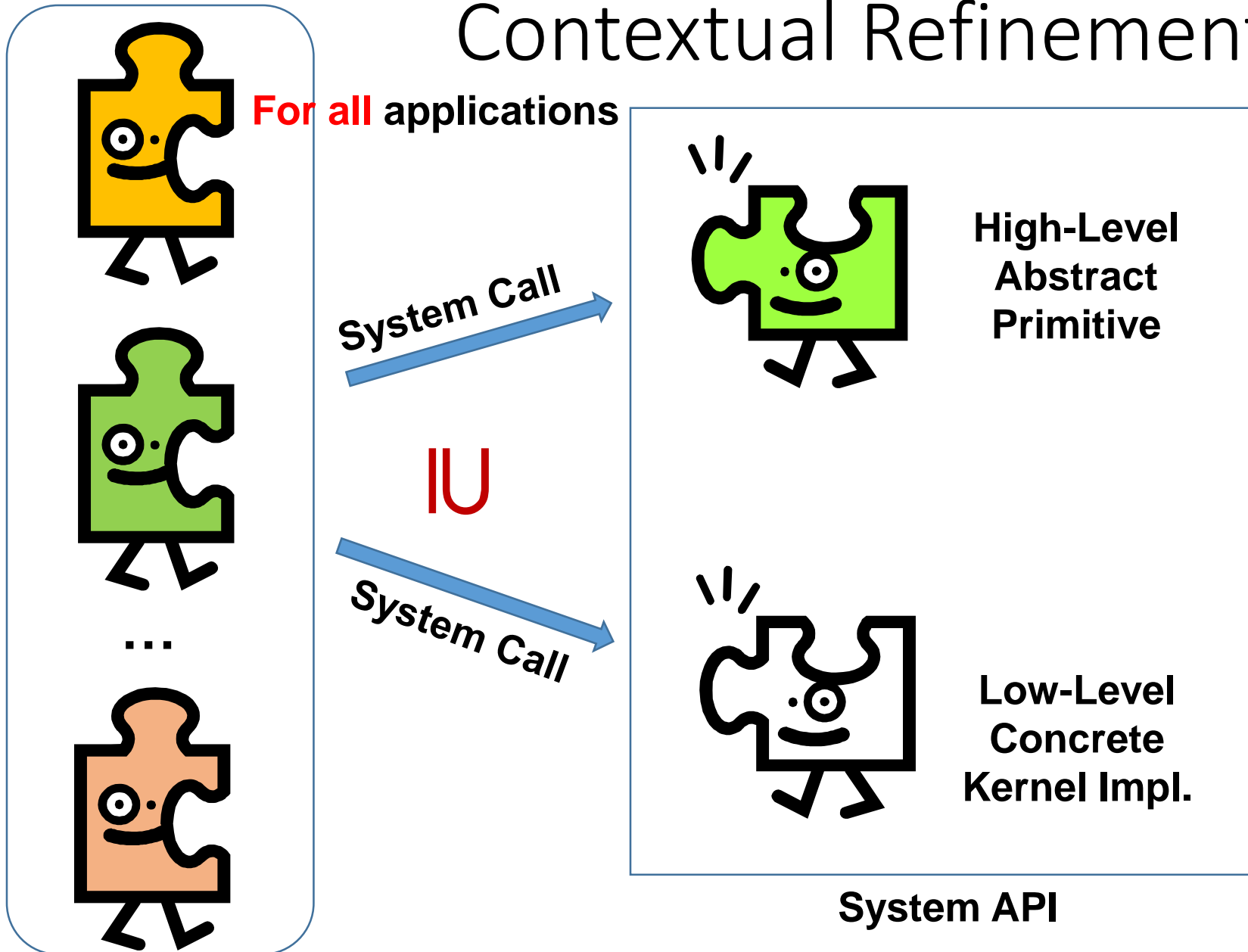
OS Correctness



Refinement



Contextual Refinement



Contextual Refinement as OS API Correctness

$O \subseteq_{\text{ctxt}} S$ iff

$$\forall A. \text{ObsBeh}(A[O]) \subseteq \text{ObsBeh}(A[S])$$

The set of observable behaviors

With some assumptions about A

A : Application O : Concrete Impl. of OS API

S : Abstract Prim.

Contextual Refinement as OS API Correctness

But OS correctness is more than API correctness:

Correctness of runtime services, e.g., scheduler
(not exported as an API)

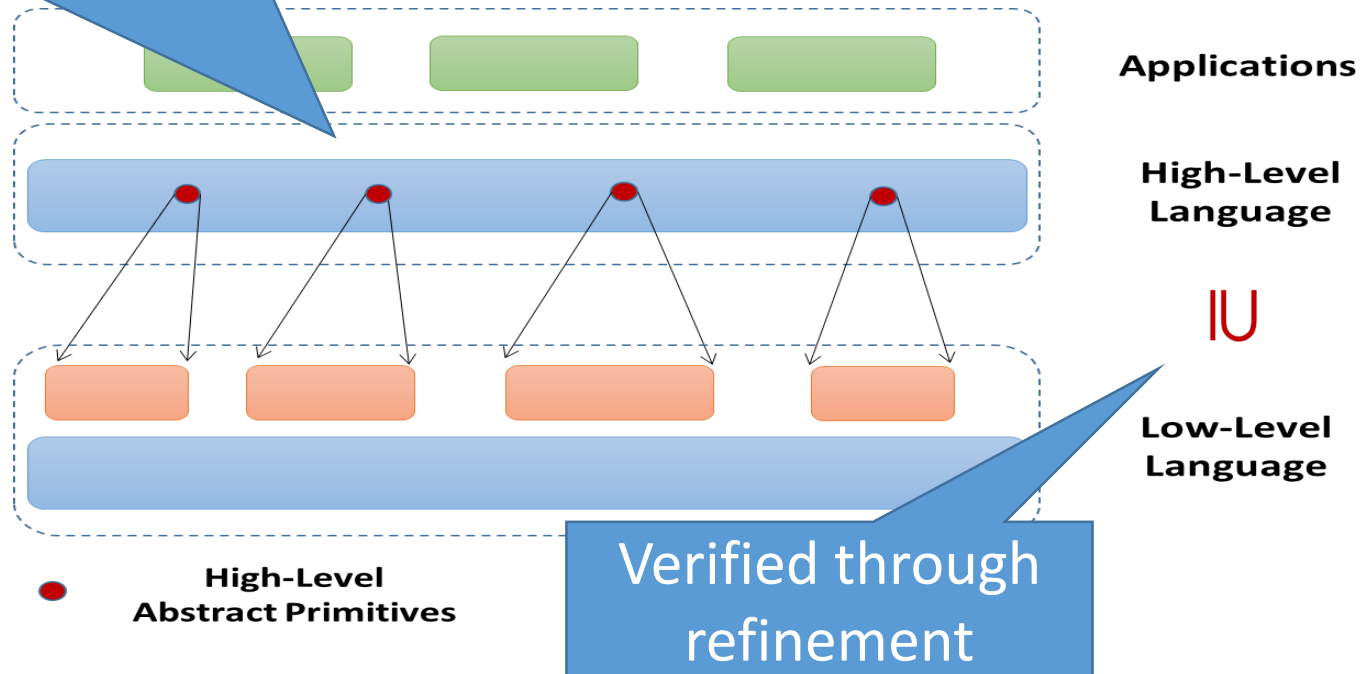
Whole system properties,
e.g., isolation and security, real-time properties, ...

Cannot be specified as abstract API primitives!

How to specify their correctness?

Runtime services and Sys. Props

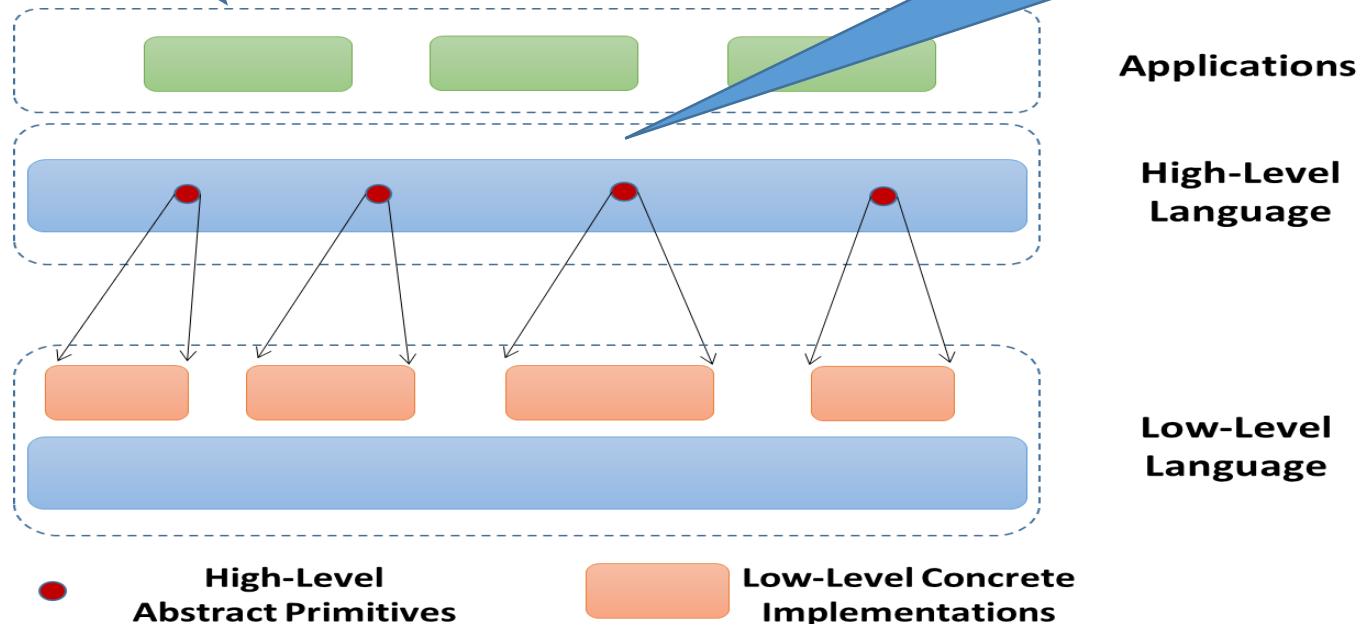
Runtime: specified as part of the high-level language semantics (e.g., scheduling)



Runtime services and Sys. Props

Whole system properties:
specified as **trace properties** of
all apps (with high-level views)

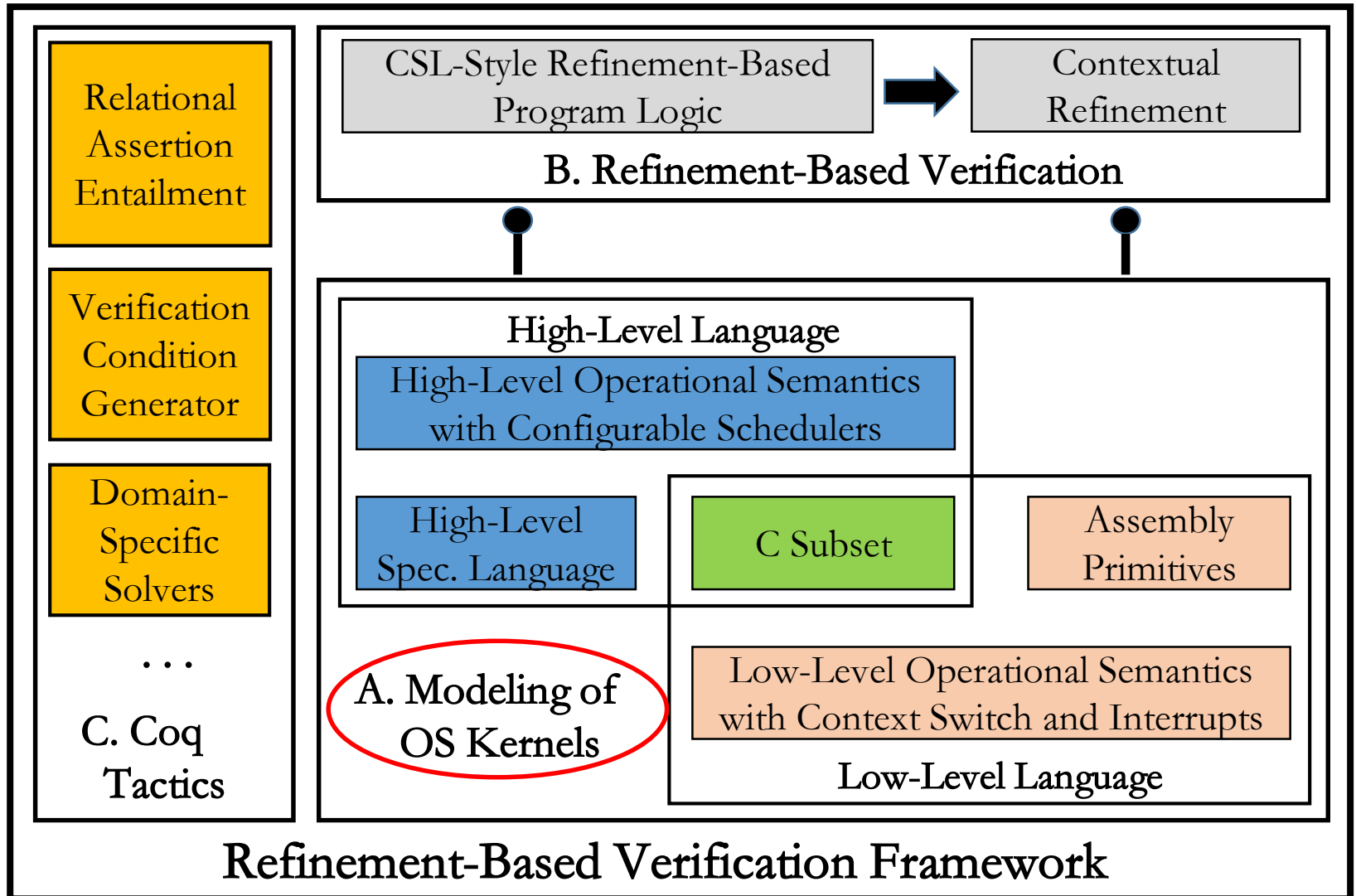
Proved at **high-level** only,
propagated to low-level
through
contextual refinement!



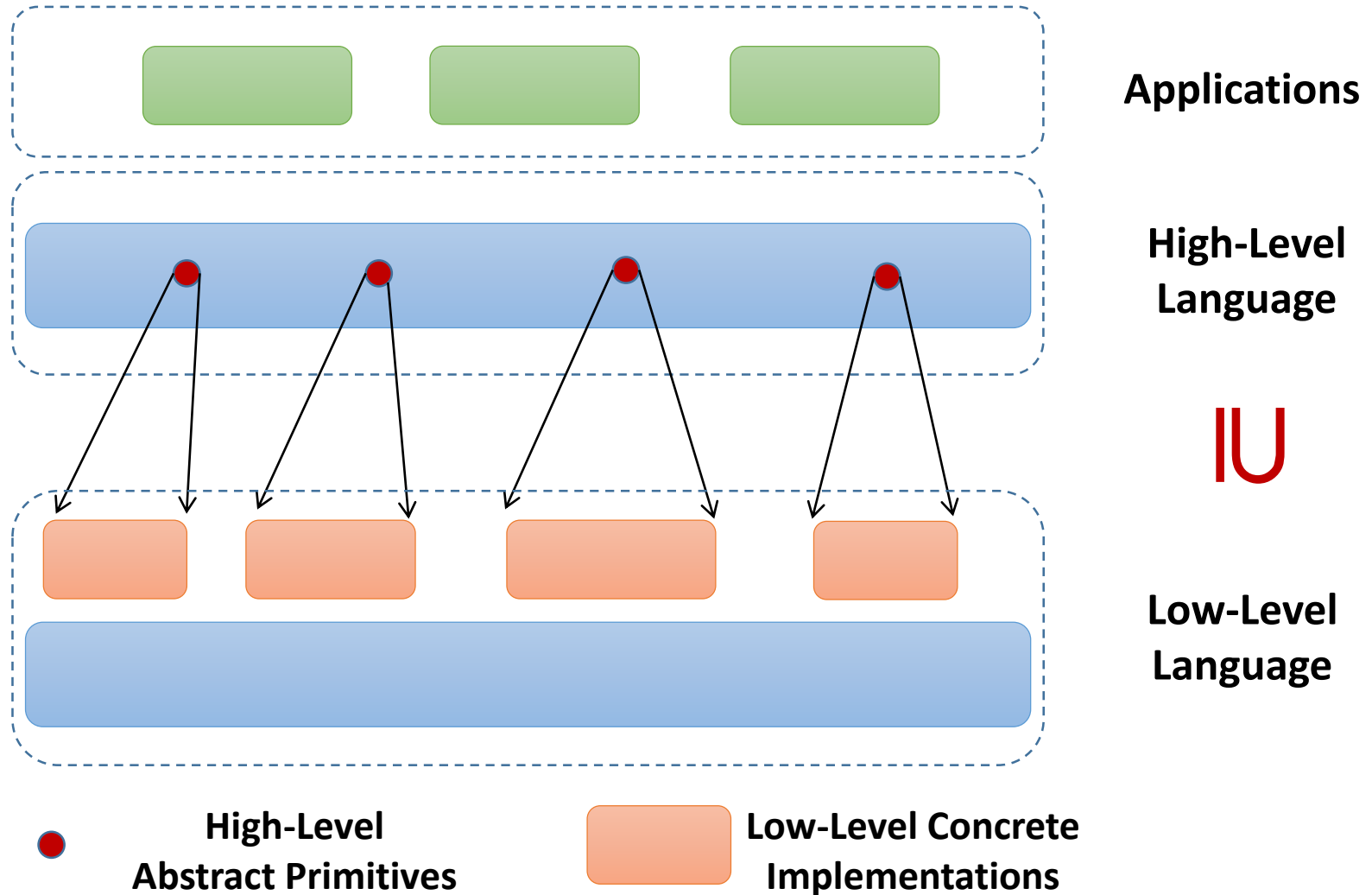
Outline

- OS Correctness Specification
- Verification Framework
 - System modeling
 - CSL-R: Program logic for refinement & multi-level interrupts
 - Coq tactics
- Verifying $\mu\text{C}/\text{OS-II}$

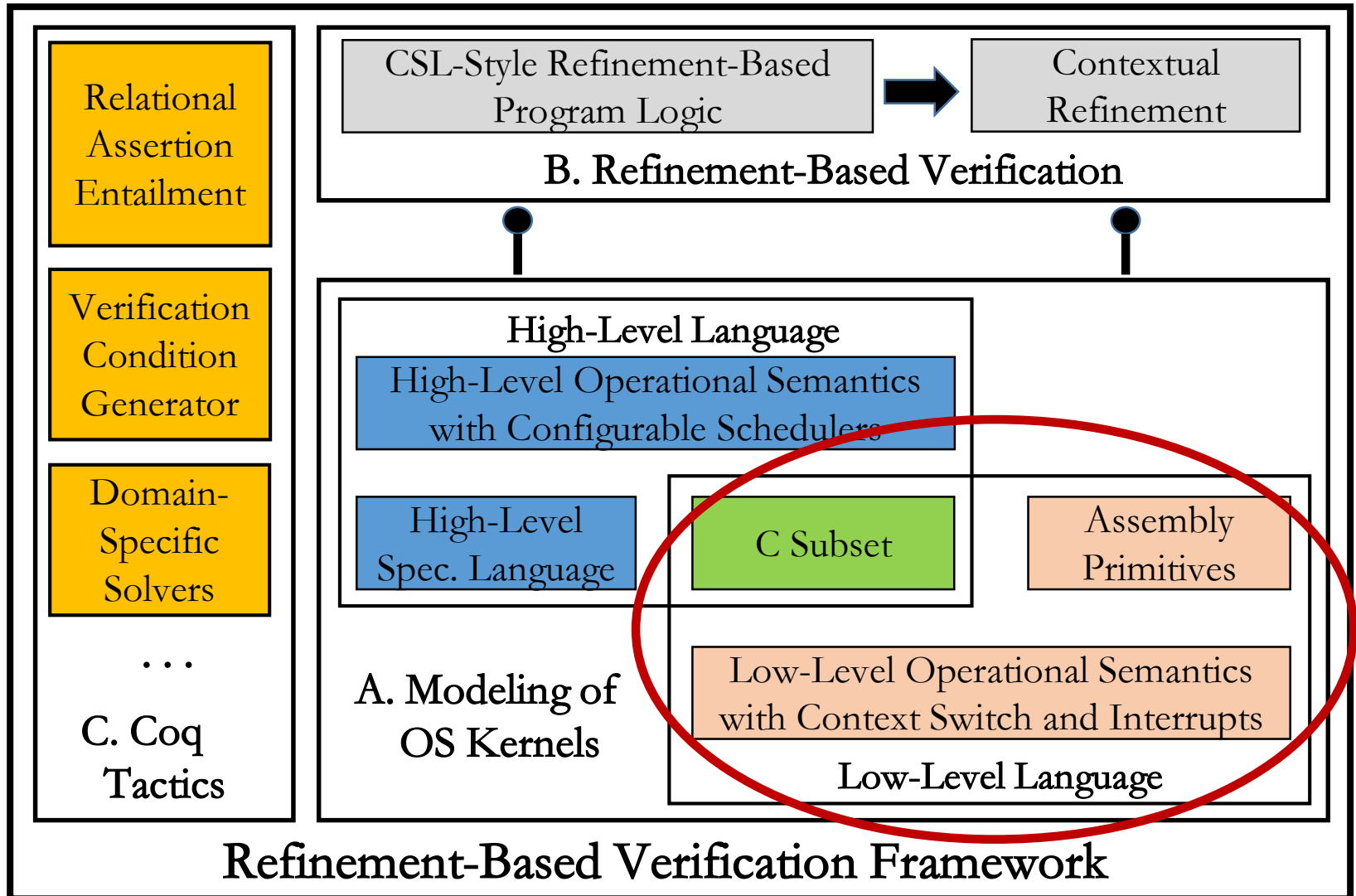
Our Verification Framework



OS Correctness



Our Verification Framework



The Low-Level Language

C Subset

$L ::= C \mid \text{Pr} \mid L;L \mid \dots$

$C ::= \text{while } e \{ C \} \mid \text{if } e \{ C1 \} \{ C2 \} \mid f(e) \mid e=e \mid \dots$

$e ::= \&e \mid *e \mid e[e] \mid e.id \mid \dots$

The Low-Level Language

```
#=====
OSCtxSw:                # Task switching from task level
    pushfl
    pushal               # Save current task's context
    mov OSTCBCur,%ebx
    mov %esp, (%ebx)     # OSTCBCur->OSTCBStkPtr = ESP
    call OSTaskSwHook    # Call user defined task switch hook
    mov OSTCBHighRdy,%eax # OSTCBCur <= OSTCBHighRdy
    mov %eax,OSTCBCur
    mov OSPrioHighRdy,%al # OSPrioCur <= OSPrioHighRdy
    mov %al,OSPrioCur
    mov OSTCBHighRdy,%ebx # ESP = OSTCBHighRdy->OSTCBStkPtr
    mov (%ebx),%esp
    popal
    popfl
    ret                  # Return to new task
#=====
```

Pr ::= encrt | excrt | switch | ...

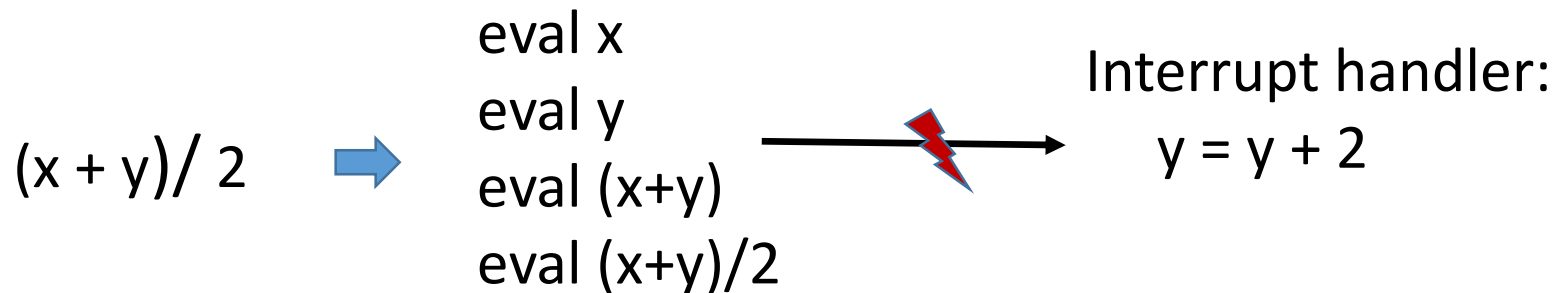
Explicit interrupts management and context switch

```
#define OS_ENTER_CRITICAL() __asm__ ("pushf \n\tcli") /* Disable interrupts*/
#define OS_EXIT_CRITICAL() __asm__ ("popf")          /* Enable interrupts*/
```

Semantics

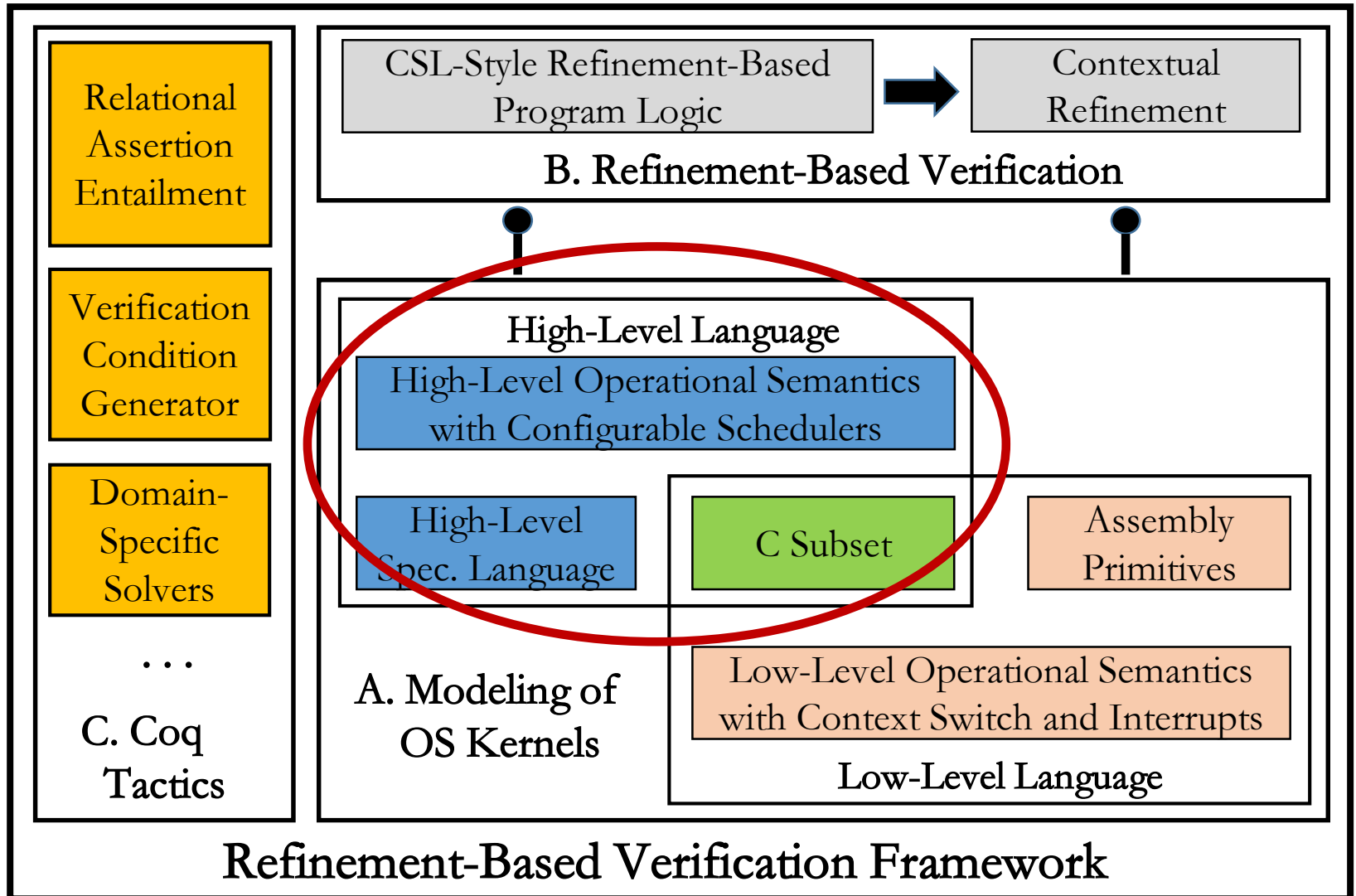
Small-step, even for expressions:

Try to be faithful to the granularity of machine-code



Semantics similar to CompCertTSO [\[Sevcik et al. 2011\]](#)
(but is interleaving semantics instead of TSO model)

Our Verification Framework



High-Level Language

C subset

High-level API
specification language

$$H ::= C \mid S \mid H;H \mid \dots$$
$$S ::= \text{sched} \mid \gamma(v) \mid S;S \mid S+S \mid \dots$$

explicit
scheduling points

High-Level Language

$$H ::= C \mid S \mid H;H \mid \dots$$

$$S ::= \text{sched} \mid \gamma(\mathbf{v}) \mid S;S \mid S+S \mid \dots$$


abstract **atomic** transitions
(over the abstract kernel states)

Example

encrt



```
void OSTimeDly (INT16U ticks)
```

```
__asm__ ("pushf \n\t cli")  
/*Disable interrupts*/
```

```
> 0) {  
    OS_ENTER_CRITICAL();
```

```
    .....
```

```
    OS_EXIT_CRITICAL();
```

Suspend the current thread,
and remove it from the
READY thread queue

```
__asm__ ("popf")  
/*Enable interrupts*/
```

```
OS_Sched();
```

```
return;
```

call scheduler

excrt



Example

Low-level Code

VS.

High-Level Spec

```
void OSTimeDly (INT16U ticks)
{
    if (ticks > 0) {
        OS_ENTER_CRITICAL();
        .....
        OS_EXIT_CRITICAL();
    }
    OS_Sched();
    return;
}
```

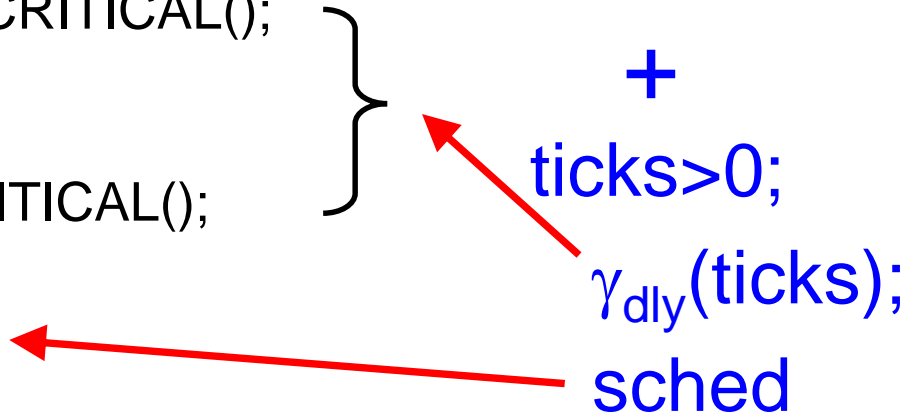
$\text{ticks} \leq 0$

+

$\text{ticks} > 0;$

$\gamma_{\text{dly}}(\text{ticks});$

sched

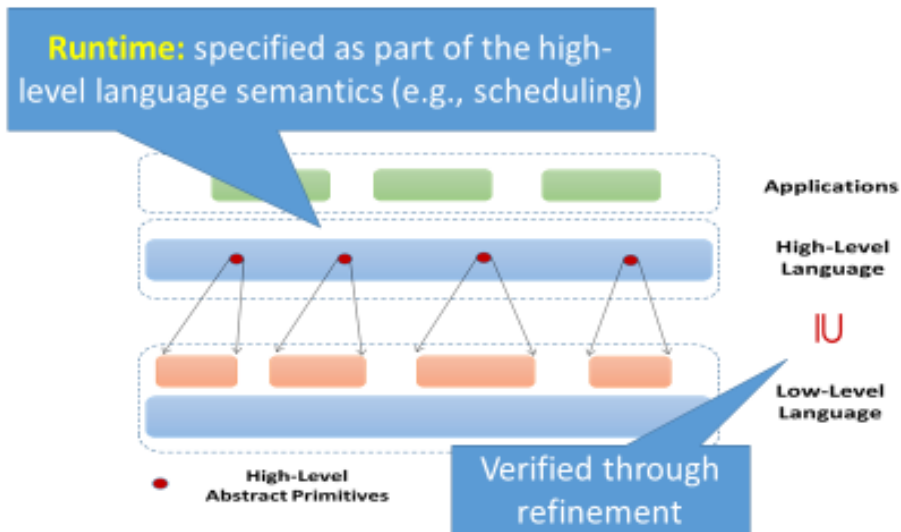


System Model

- Low-level impl. **O**: (η_a , θ , η_i)
 - η_a : API implementations
 - θ : Interrupt handlers
 - η_i : Internal functions
- High-level spec. **S**: (φ , ε , χ)
 - φ : API specs. (high-level primitives for APIs)
 - ε : Abstract events (high-level primitives for int. handlers)
 - χ : **Abstract scheduler**
 - Scheduling policy can be customized by instantiating χ

System Model

Runtime services and Sys. Props



for APIs)

imitives for int. handlers)

- χ : **Abstract scheduler**
 - Scheduling policy can be customized by instantiating χ
 - Shows abstractions for runtime

System Model

- Low-level impl. **O**: (η_a, θ, η_i)

- η_a : API implementations
- θ : Interrupt handlers
- η_i : Internal functions

Verification goal:

$$(\eta_a, \theta, \eta_i) \subseteq_{\text{ctxt}} (\varphi, \varepsilon, \chi)$$

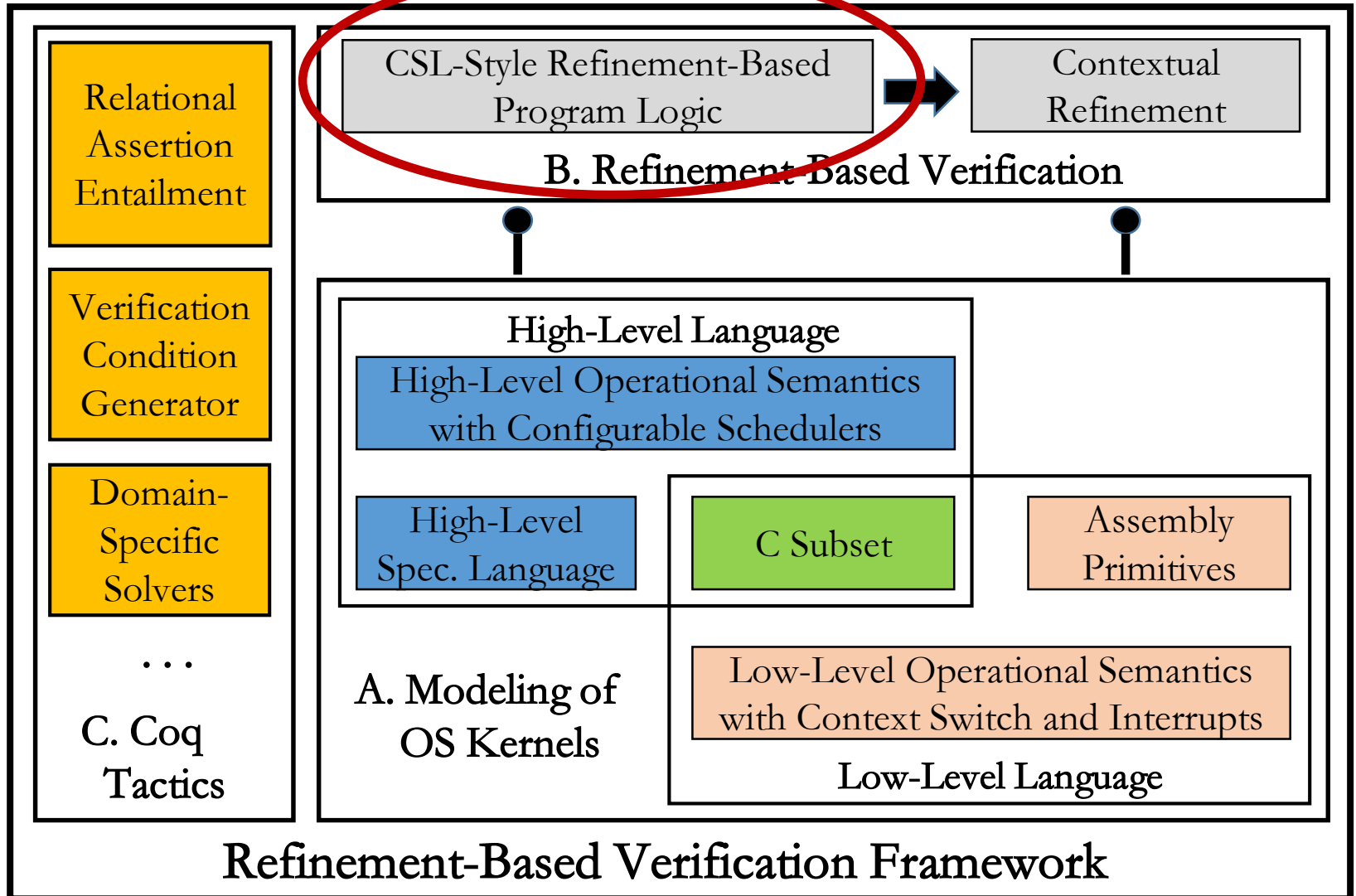
- High-level spec. **S**: $(\varphi, \varepsilon, \chi)$

- φ : API specs. (high-level primitives for APIs)
- ε : Abstract events (high-level primitives for int. handlers)
- χ : Abstract scheduler
 - Scheduling policy can be customized by instantiating χ
 - Shows abstractions for runtime

Outline

- OS Correctness Specification
- Verification Framework
 - System modeling
 - CSL-R: Program logic for refinement & multi-level interrupts
 - Coq tactics
- Verifying $\mu\text{C}/\text{OS-II}$

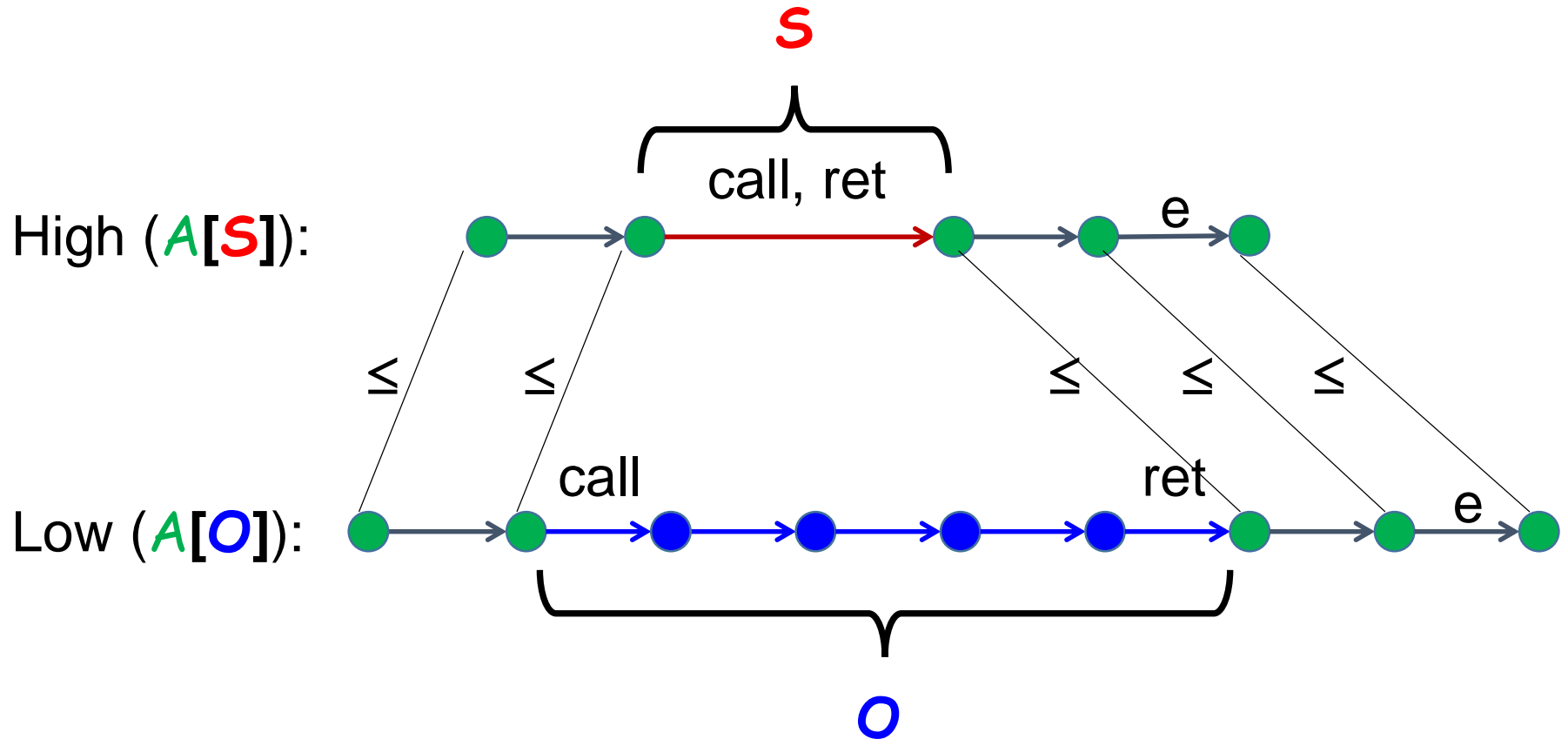
Our Verification Framework



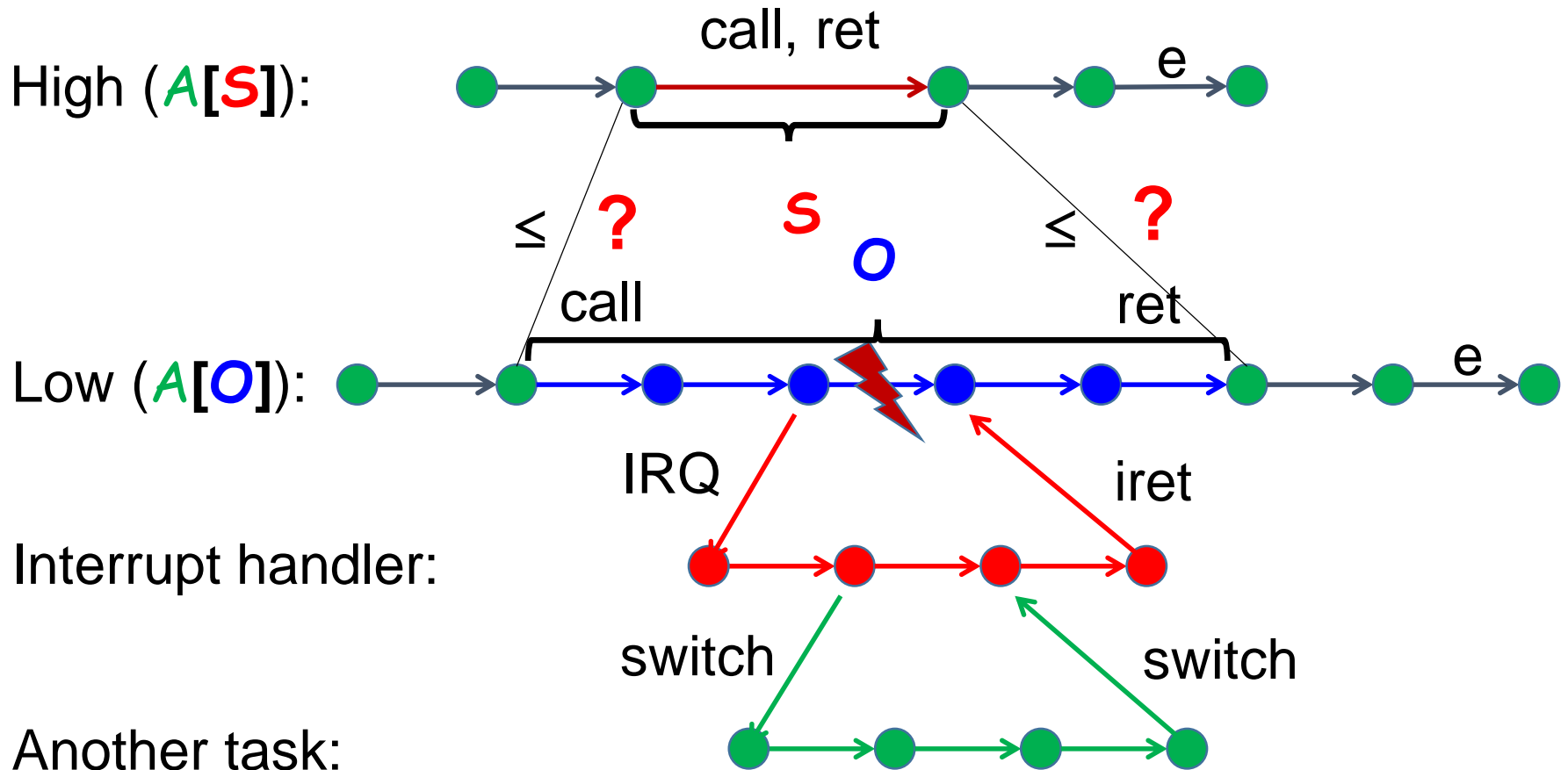
Program Logic for Refinement and Multi-Level Interrupts

- Relational program logic for simulation/refinements
[Liang et al. PLDI'13, CSL-LICS'14]
- Ownership-Transfer semantics for interrupts
[Feng et al. PLDI'08]
- Combining the two: CSL-R for refinement reasoning with multi-level interrupts

Refinement Verification via Simulation

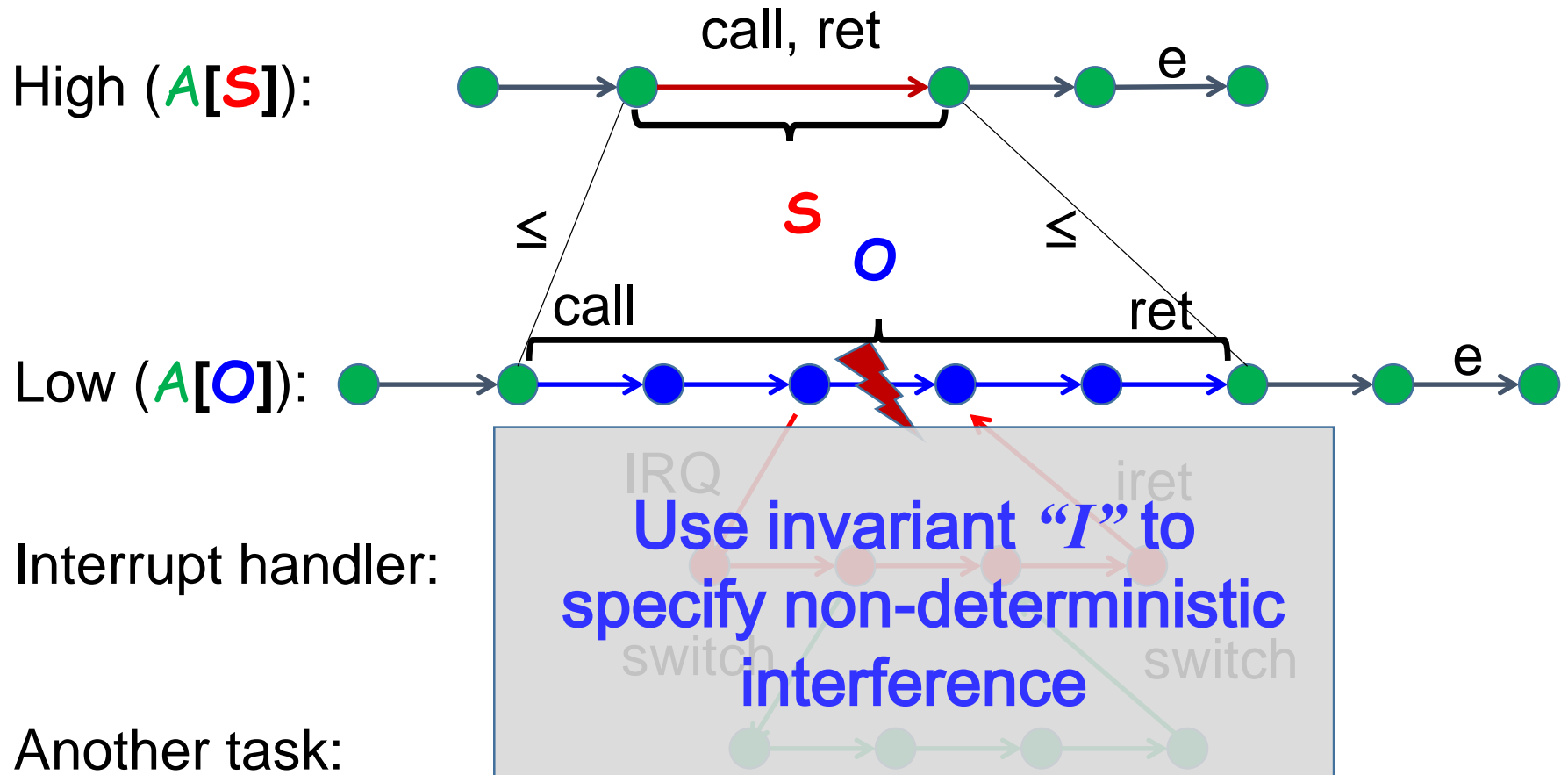


Simulation with Interrupts & Multitasking

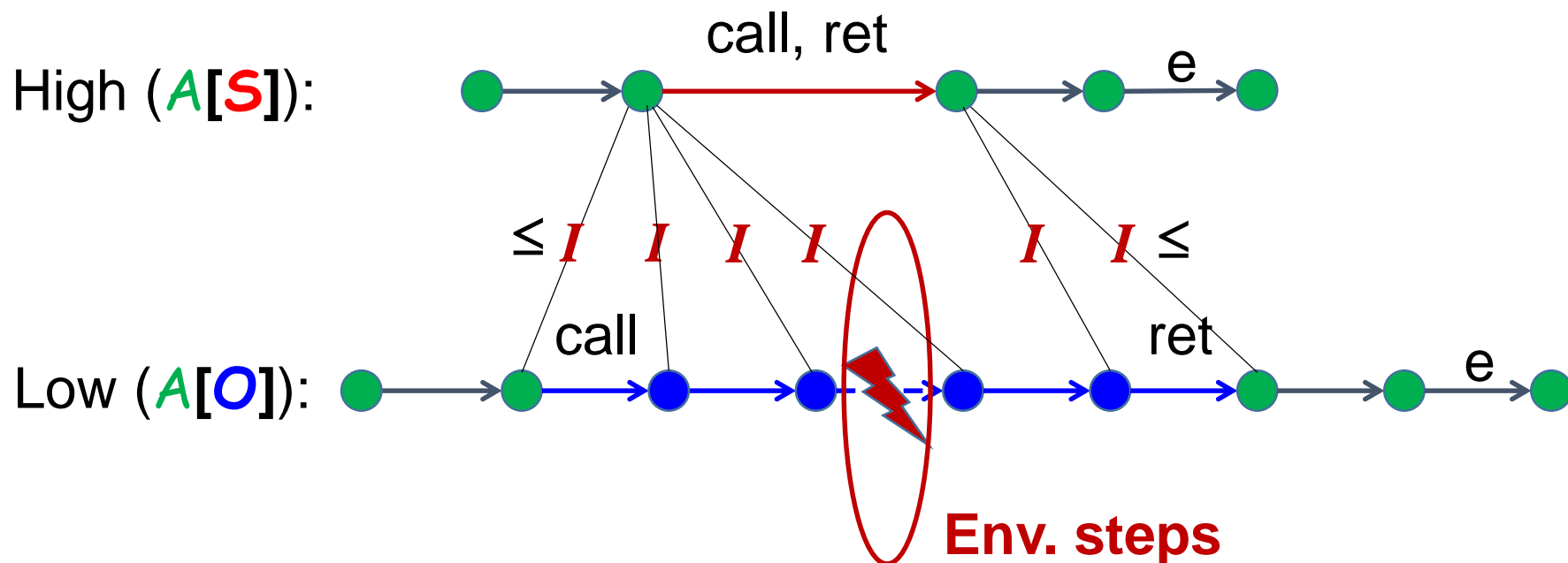


How to do compositional verification?

Simulation with Interrupts & Multitasking

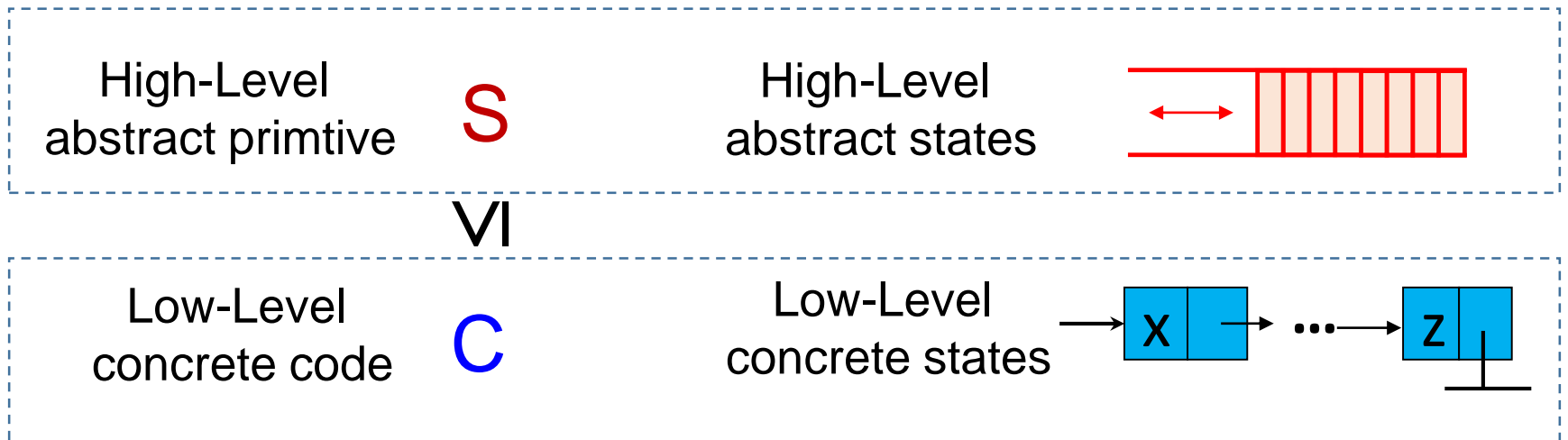


Simulation with Interrupts & Multitasking



Adapted from RGSim [Liang et al. POPL'12] and
the relational program logic [Liang et al. PLDI'13, CSL-LICS'14]

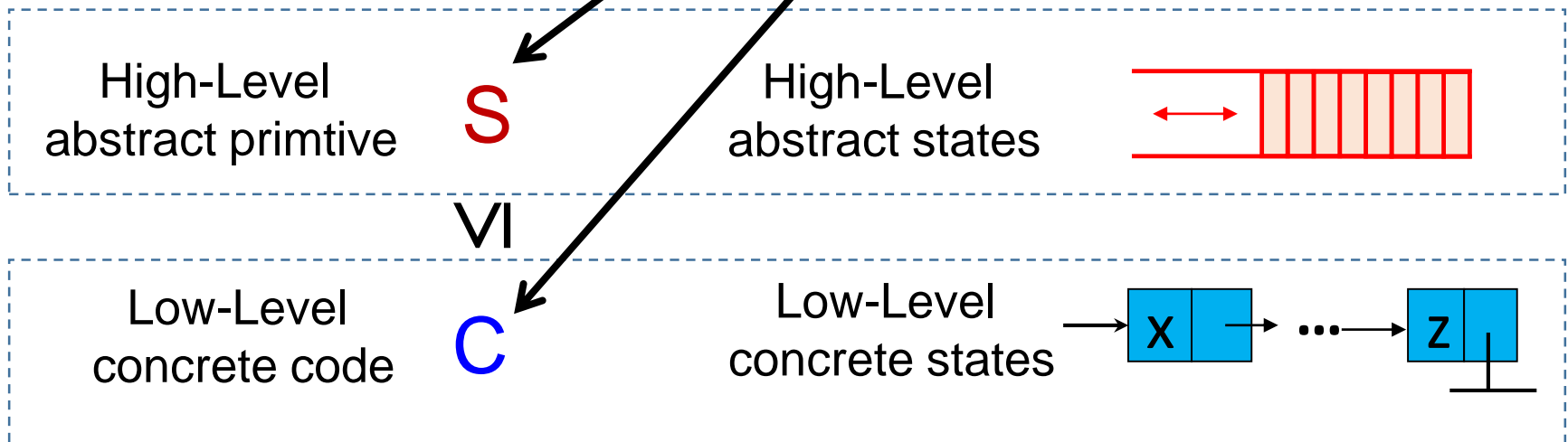
Program Logic for Simulation



Program Logic for Simulation

- Judgement

Remaining high-level code
that needs to be refined

$$I \vdash \{p * \mathbf{[|S|]}\} \mathbf{C} \{q * \mathbf{[|end|]}\}$$


Program Logic for Simulation

- Judgement

No remaining high-level code (refinement is done)

$$I \vdash \{p * \mathbf{[|S|]}\} \mathbf{C} \{q * \mathbf{[|end|]}\}$$

High-Level
abstract primitive

S

High-Level
abstract states

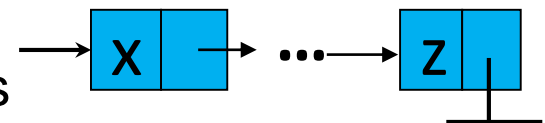


V

Low-Level
concrete code

C

Low-Level
concrete states



Program Logic for Simulation

- Judgement

Relational assertions for
pre-/post- condition

$$I \vdash \{p^* \text{[|S|]} \} C \{q^* \text{[|end|]} \}$$

High-Level
abstract primitive

S

High-Level
abstract states

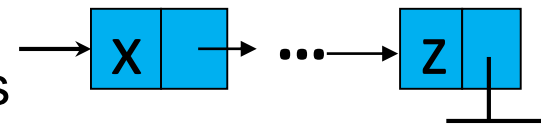


V

Low-Level
concrete code

C

Low-Level
concrete states

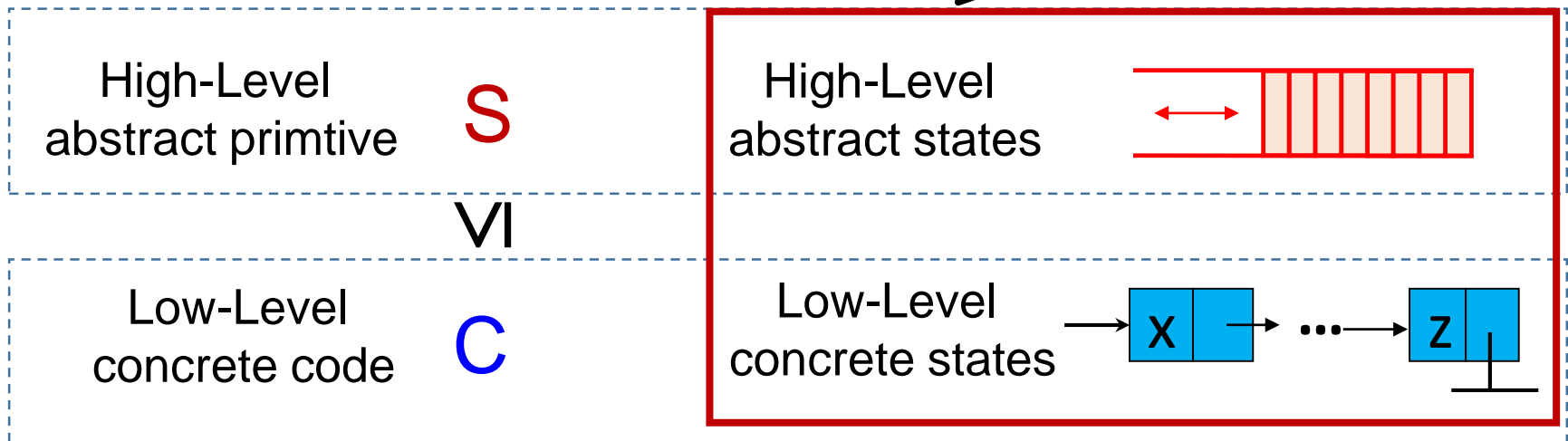


Program Logic for Simulation

- Judgement

Relational Invariants

$$I \vdash \{p * \mathbf{[|S|]}\} \text{ } \mathbf{C} \{q * \mathbf{[|end|]}\}$$



Soundness

If

$$I \vdash \{p^* [|\mathbf{S}|]\} \mathbf{C} \{q^* [|\mathbf{end}|]\}$$

then \mathbf{C} is simulated by \mathbf{S} , ...

An Example

```
void Add() {
```

```
    OS_ENTER_CRITICAL();
```

```
    Count ++ ;
```

```
    OS_EXIT_CRITICAL();
```

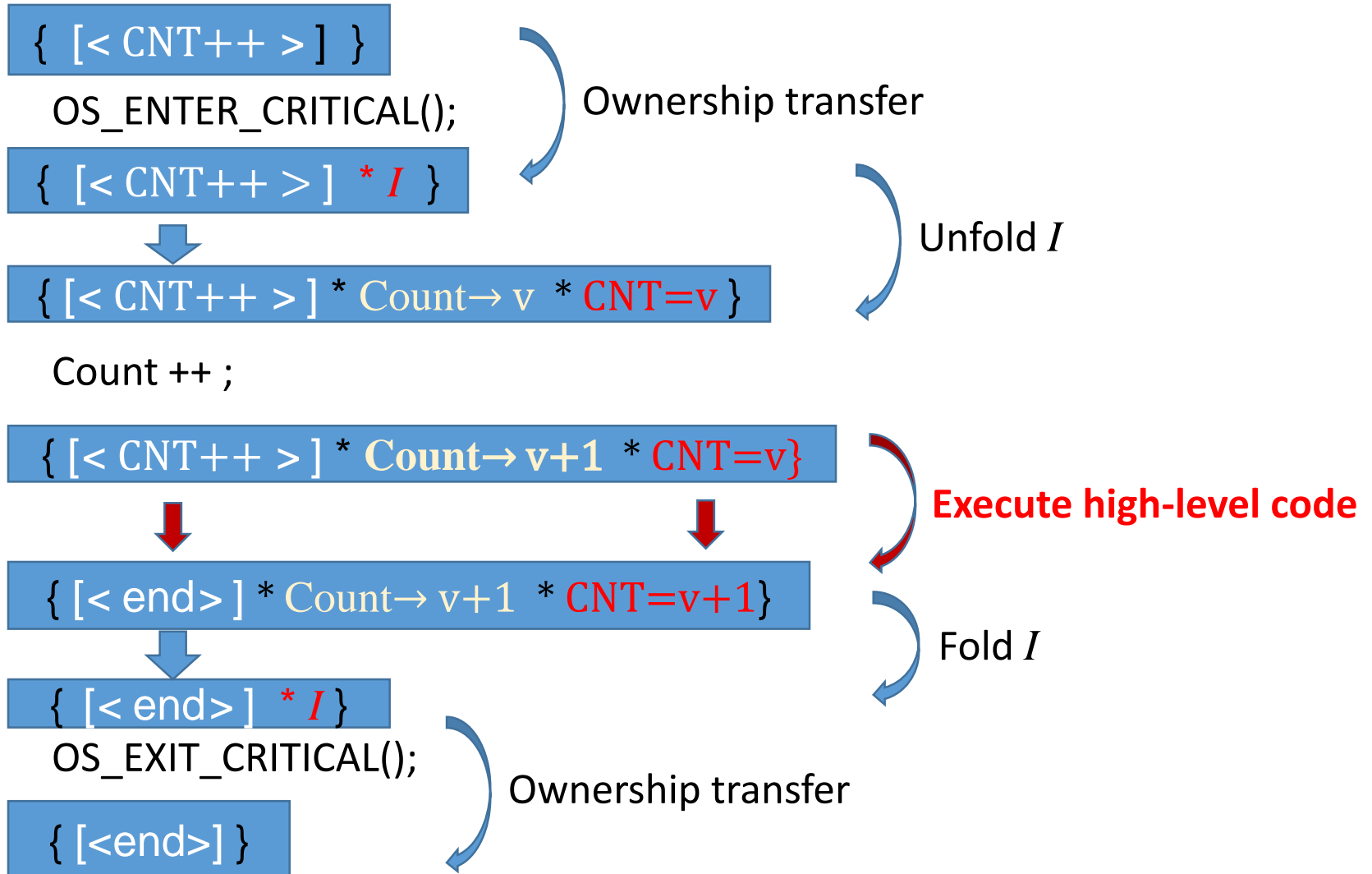
```
}
```

\leq

$\langle \text{CNT}++ \rangle$

$I ::= \exists v. \text{Count} \rightarrow v * \text{CNT} = v$

An Example



An Example

```
{ [<CNT++>] }
```

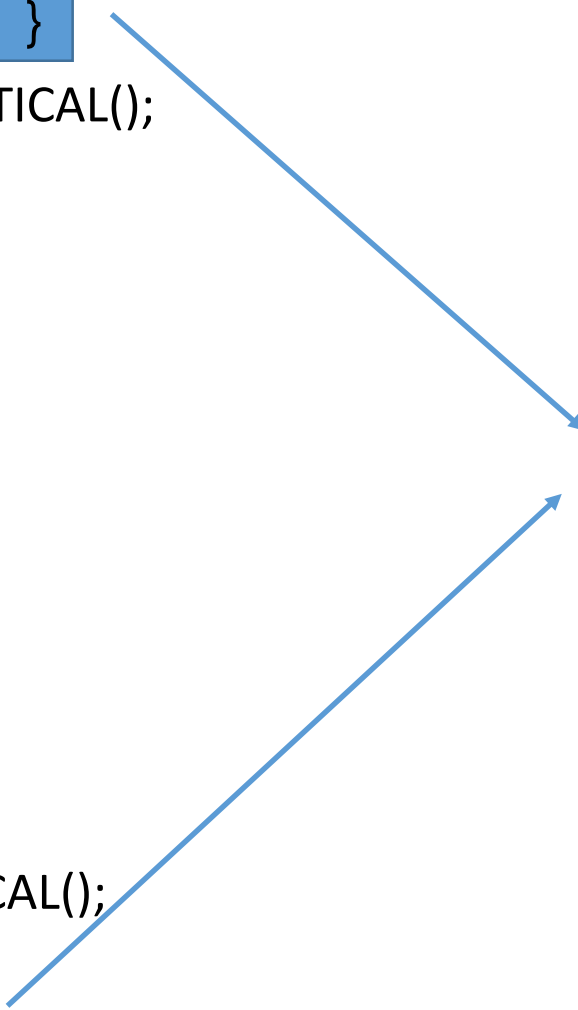
```
OS_ENTER_CRITICAL();
```

```
Count ++ ;
```

```
OS_EXIT_CRITICAL();
```

```
{ [<end>] }
```

The code refines
<CNT++>



An Example

```
{ [< CNT++ > ] }
```

```
OS_ENTER_CRITICAL();
```

```
{ [< CNT++ > ] * Count→ v * CNT=v }
```

```
Count ++ ;
```

```
{ [< CNT++ > ] * Count→ v+1 * CNT=v }
```



```
{ [< end> ] * Count→ v+1 * CNT=v+1 }
```



Execute high-level code

```
OS_EXIT_CRITICAL();
```

```
{ [<end> ] }
```

An Example

{ [< CNT++ >] }

OS_ENTER_CRITICAL();

Count ++ ;

{ [< CNT++ >] * Count → v+1 * CNT=v }



{ [< end >] * Count → v+1 * CNT=v+1 }



Execute high-level code

OS_EXIT_CRITICAL();

{ [< end >] }

Abstract consequence rule:

$$\frac{p^*[S] \Rightarrow r^*[S'] \quad \vdash \{ r^*[S'] \} C \{ q \}}{\vdash \{ p^*[S] \} C \{ q \}}$$

$$p \Rightarrow q \text{ iff } \forall (\sigma, \Sigma, S) \models p, \\ \exists (\Sigma', S'). (\Sigma, S) \rightarrow^* (\Sigma', S') \\ \wedge (\sigma, \Sigma', S') \models q,$$

An Example

```
{ [< CNT++ > ] }
```

```
OS_ENTER_CRITICAL();
```

```
{ [< CNT++ > ] * I }
```



Ownership transfer

```
Count ++ ;
```

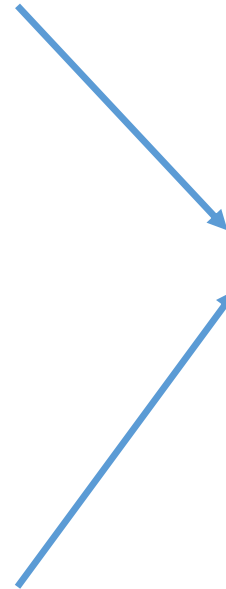
```
{ [< end> ] * I }
```

```
OS_EXIT_CRITICAL();
```

```
{ [<end> ] }
```



Ownership transfer



Interrupt reasoning

Interrupt Reasoning

Program invariant [O'Hearn CONCUR'04]

There is always a **partition** of resource among concurrent entities, and each concurrent entity only accesses its own part.



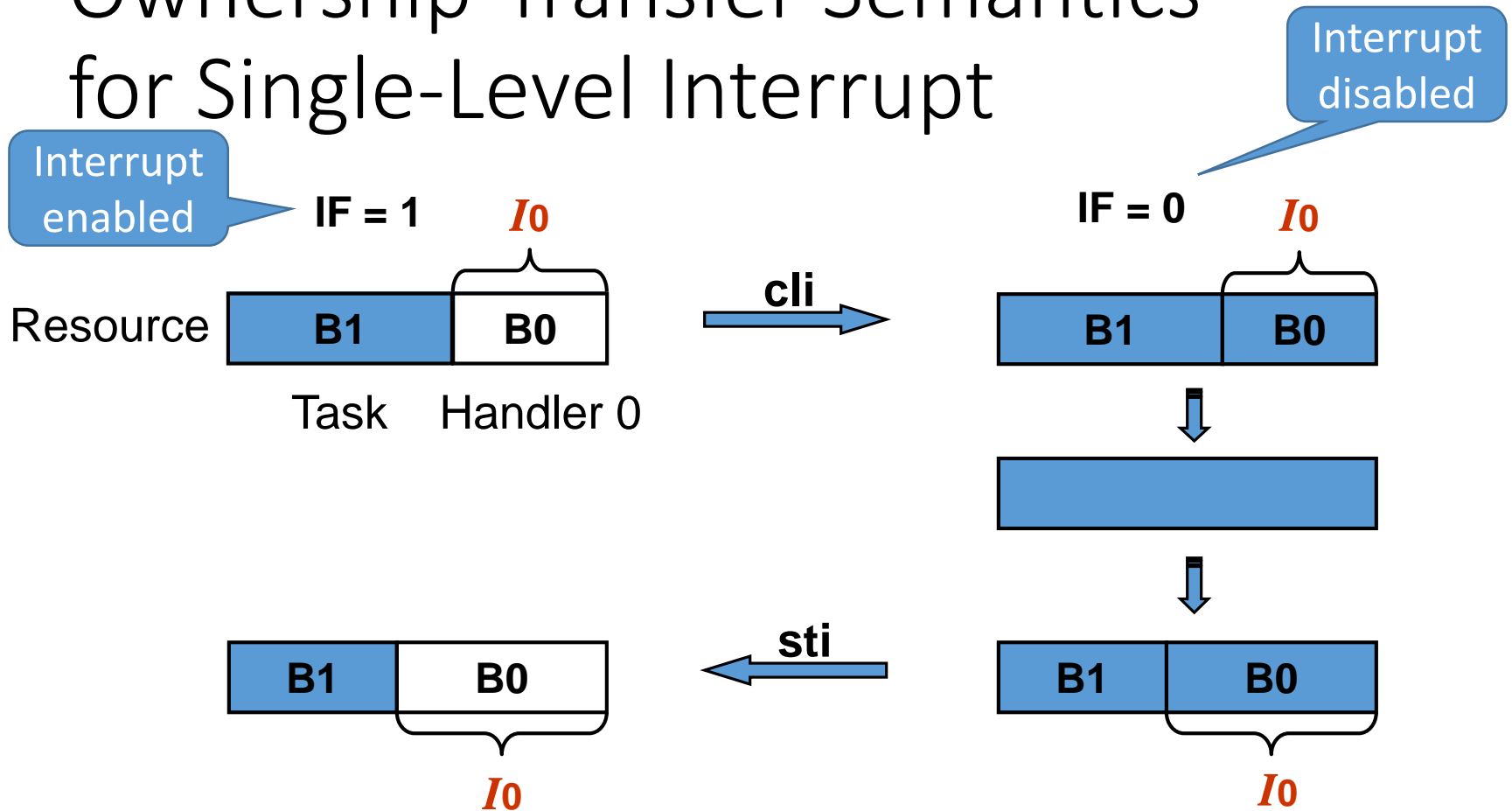
Tasks and interrupt
handlers

But note:

The partition is dynamic: ownership of resource can be dynamically transferred.

Interrupt operations can be modeled as operations that trigger resource ownership transfer. [Feng et al. PLDI'08]

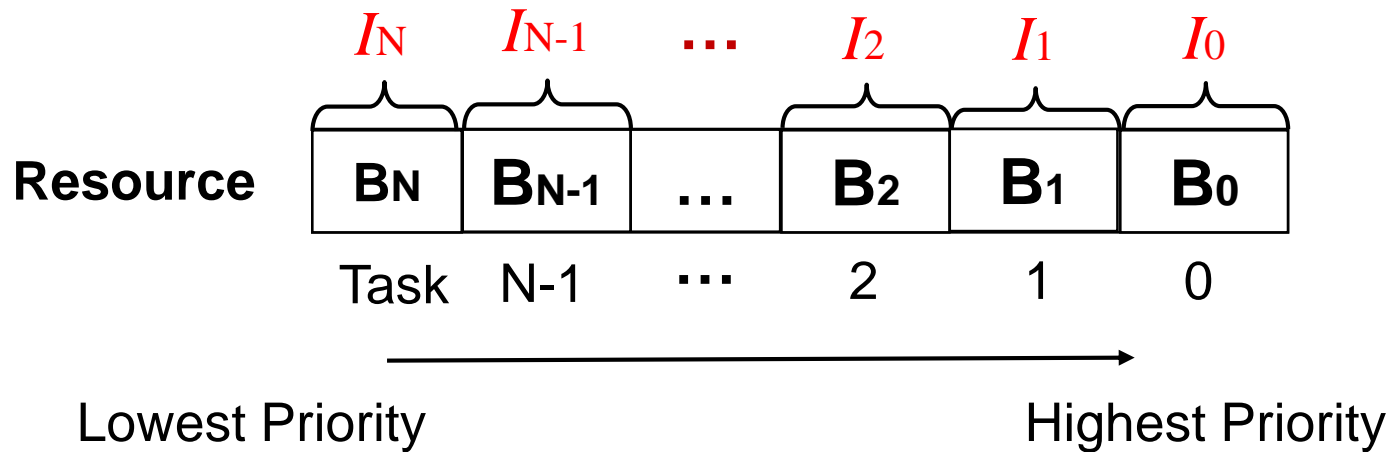
Ownership-Transfer Semantics for Single-Level Interrupt



$I_0 \vdash \{p\} \text{ cli } \{p * I_0\}$

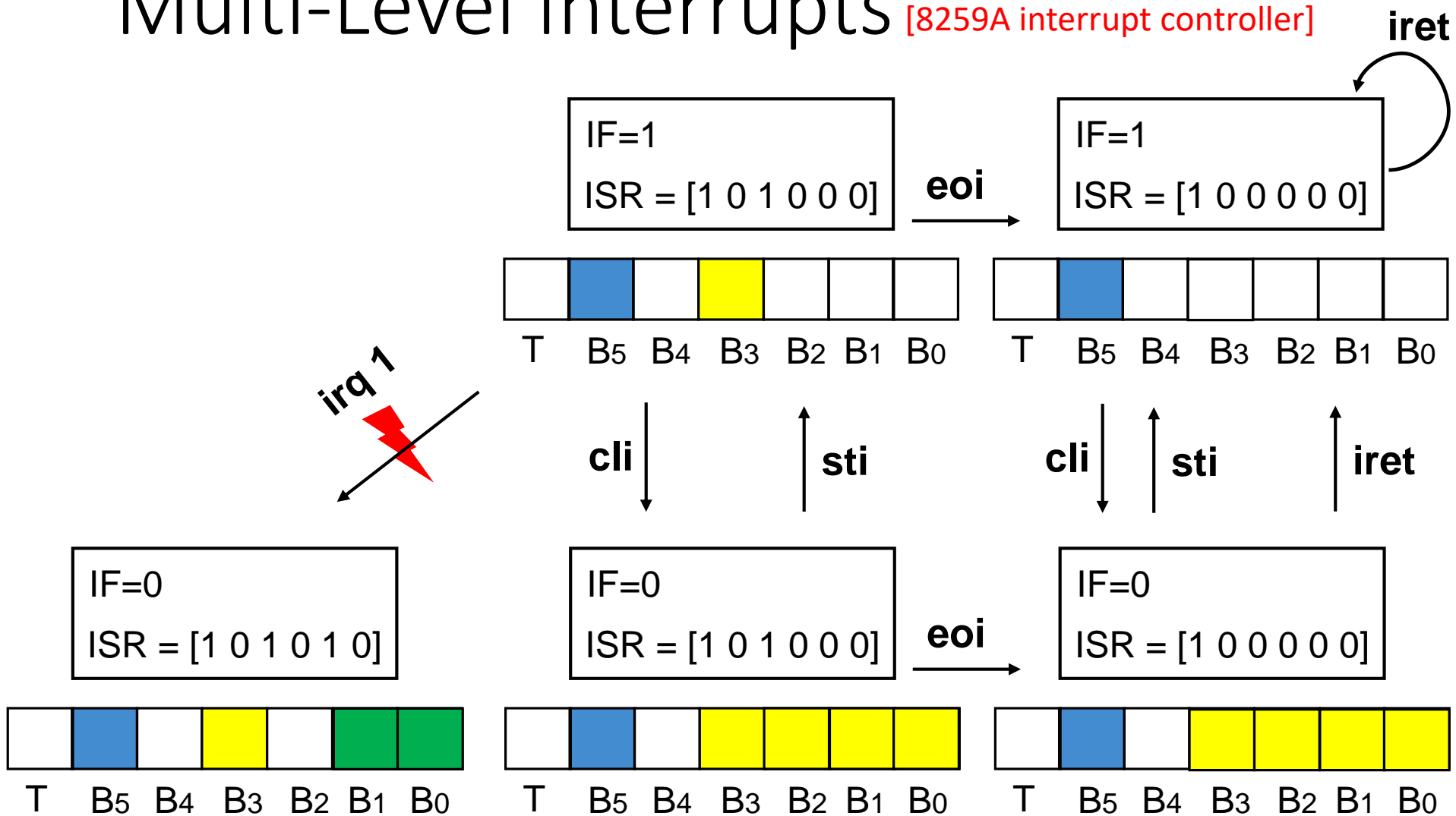
$I_0 \vdash \{p * I_0\} \text{ sti } \{p\}$

Memory Model for Multi-Level Interrupts

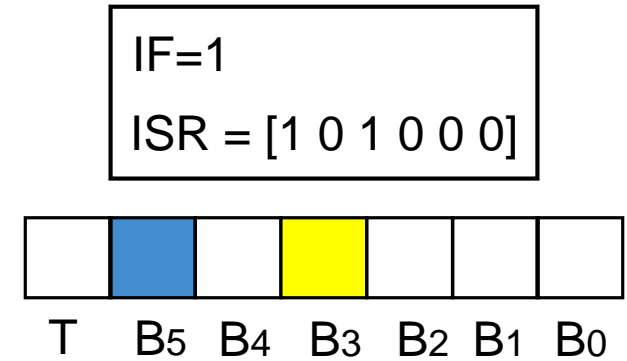


- Higher-priority handler has priority to select its required resource
- N blocks are assigned to N interrupt handlers
- Each well-formed resource block is specified by a resource invariant

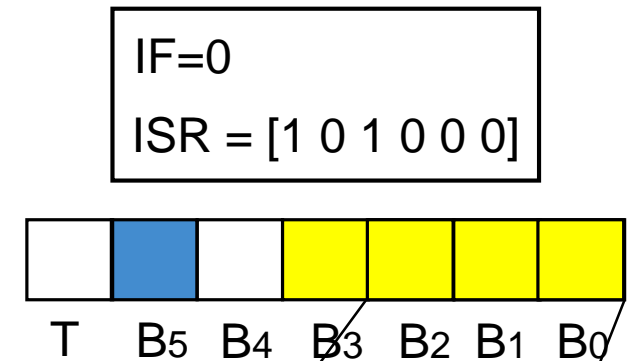
Ownership-Transfer Semantics for Multi-Level Interrupts [8259A interrupt controller]



Inference Rules for Interrupt Operations



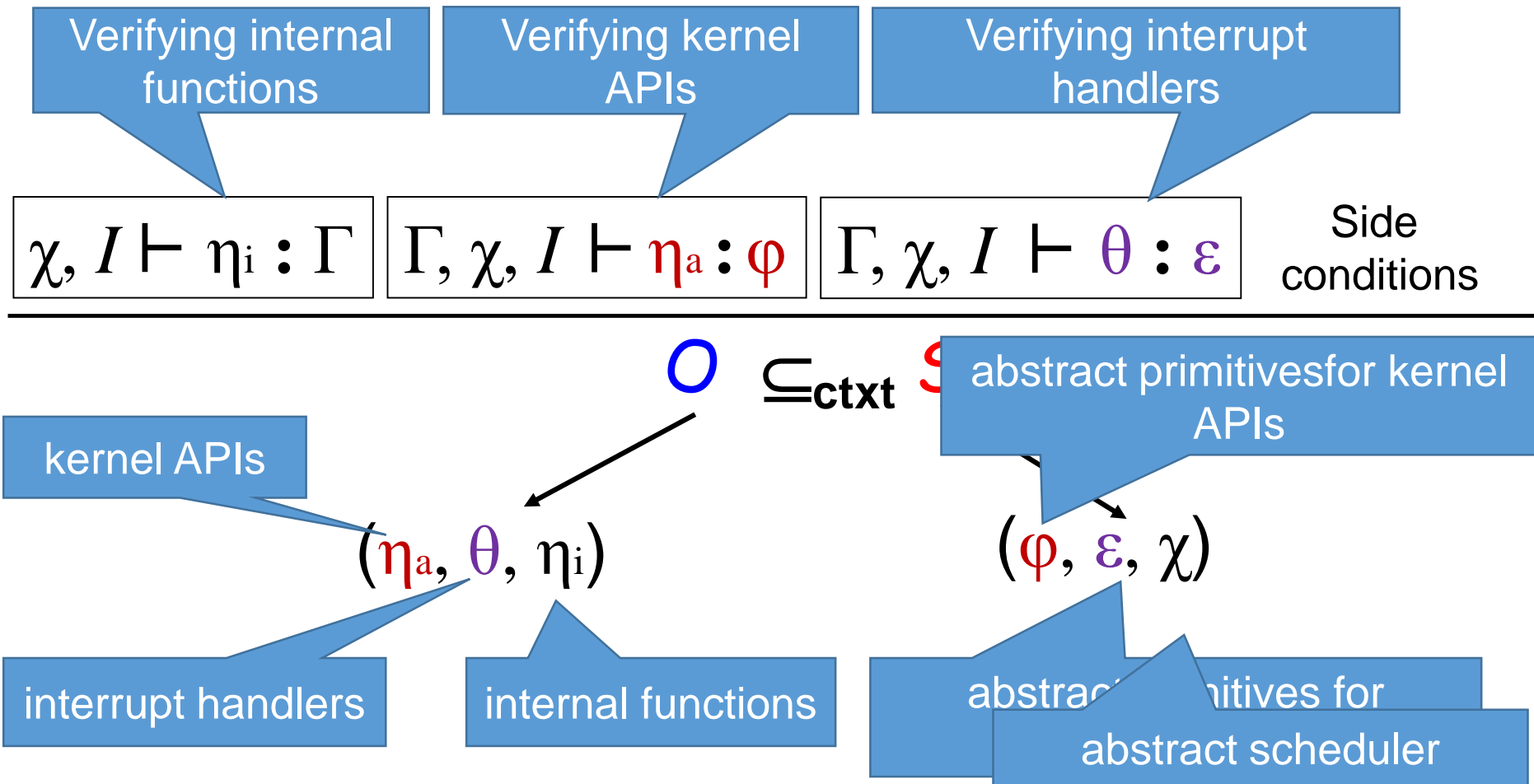
cli
↓



ISR(k) = 1

$I \vdash \{ [ISR, 1, k]^* p \} \text{ cli } \{ [ISR, 0, k]^* p^* \text{ } \boxed{I[0 \dots k-1]} \}$

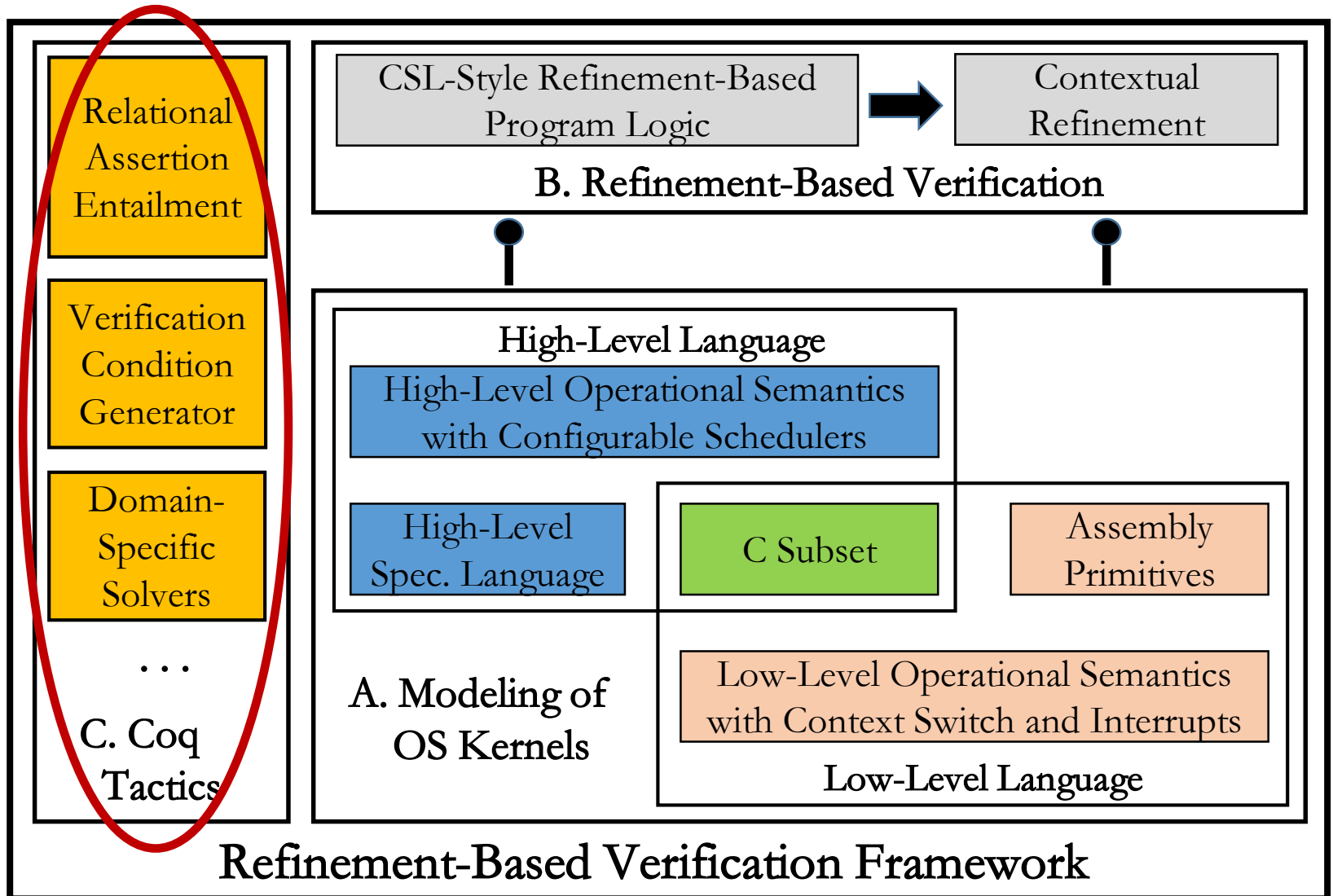
Top Rule for Proving $O \sqsubseteq_{\text{ctxt}} S$



Outline

- OS Correctness Specification
- Verification Framework
 - System modeling
 - CSL-R: Program logic for refinement & multi-level interrupts
 - Coq tactics
- Verifying $\mu\text{C}/\text{OS-II}$

Our Verification Framework



Coq Tactics for Automation Support

- Verification condition generator : **hoare forward**
 - Automatically select and apply the inference rules
- Assertion entailment prover : **sep auto**
 - Automatically prove “ $p \Rightarrow q$ ”
- Domain specific solvers : **mauto ...**

$$\forall x. x < 64 \rightarrow x \gg 3 < 8 ; \quad \forall x. x < 8 \rightarrow (x \ll 3) \& 7 = 0$$

Coq Tactics for Automation Support

- To reduce the proof efforts
- To hide the underlying details of the verification framework
- To prove domain specific propositions

The ratio of Coq scripts to the verified C is around 27:1

lots of space for improvement

Outline

- OS Correctness Specification
- Verification Framework
 - System modeling
 - CSL-R: Program logic for refinement & multi-level interrupts
 - Coq tactics
- Verifying $\mu\text{C}/\text{OS-II}$

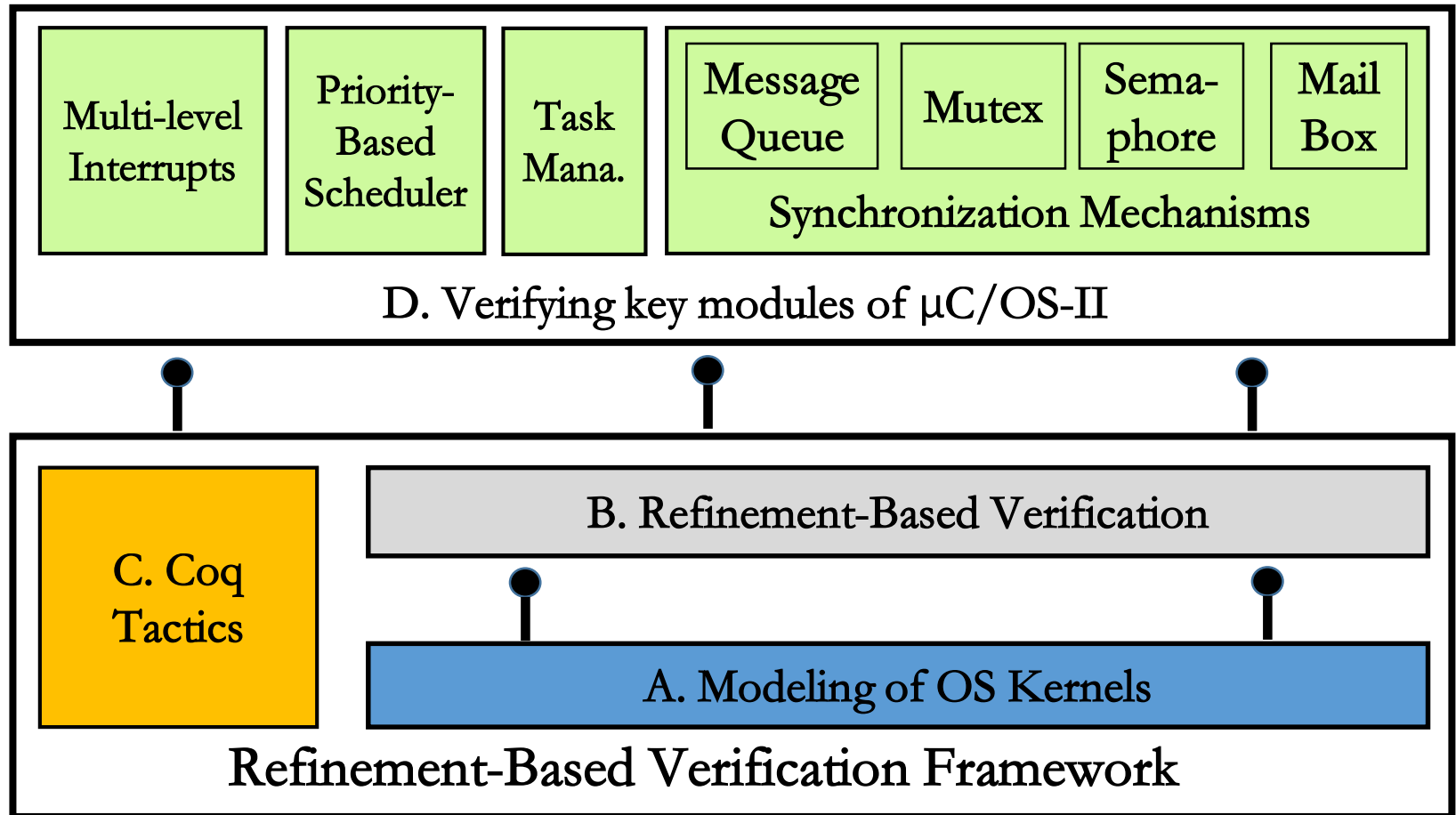
μC/OS-II

- A commercial preemptive real-time multitasking OS kernel developed by Micrium.

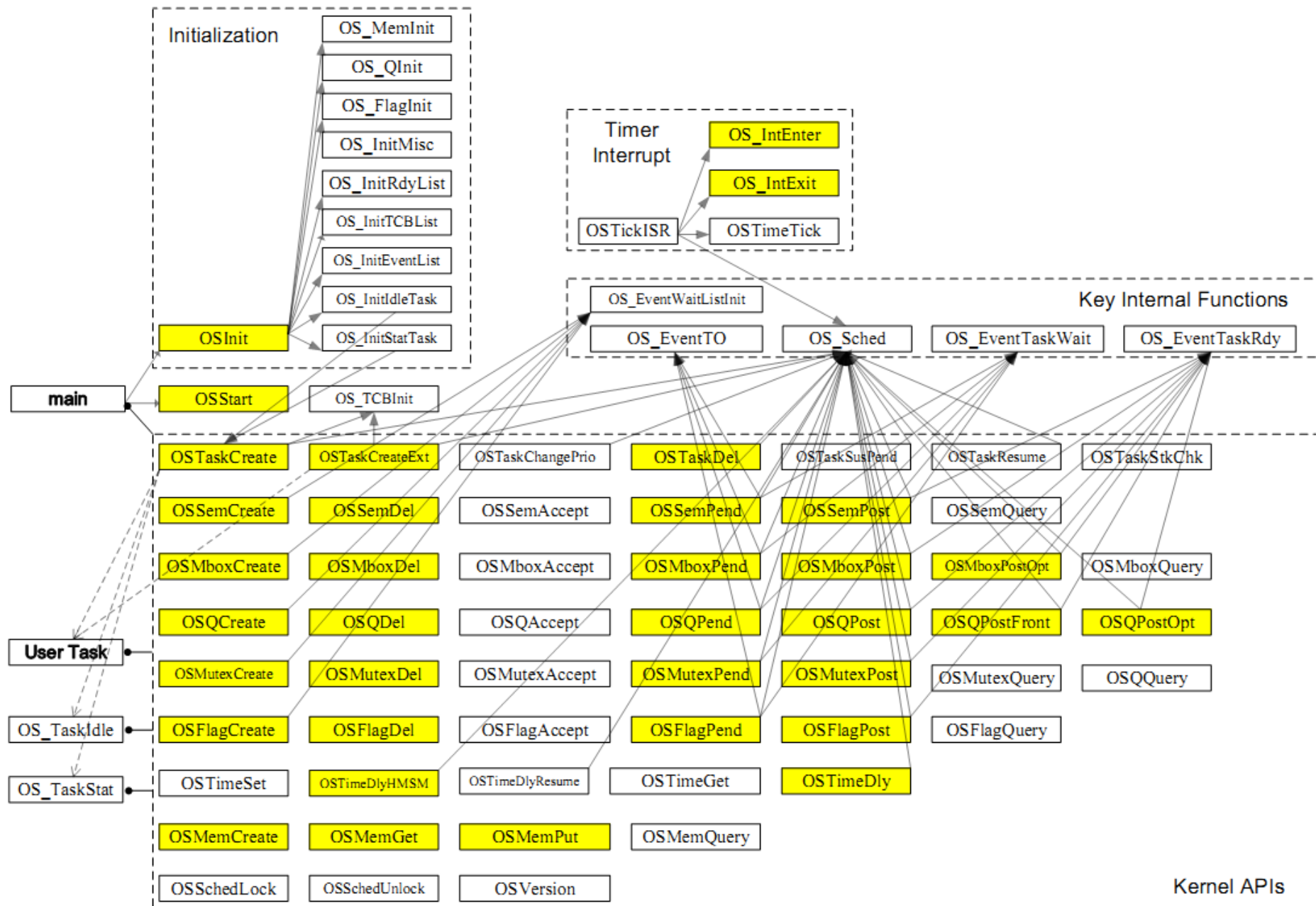


- 6,316 lines of C & 316 lines of assembly code.
- Multitasking & Multi-Level interrupts & Preemptive priority-based scheduling & Synchronization mechanism
- Deployed in many real-world safety critical applications
 - Avionics and medical equipments, *etc.*

Verifying $\mu\text{C}/\text{OS-II}$



Frequently Used APIs :



Verified APIs : [REDACTED]

covers 68% of the

Timer interrupt handler APIs

Scheduler

Task management

Key Internal Functions

Semaphore

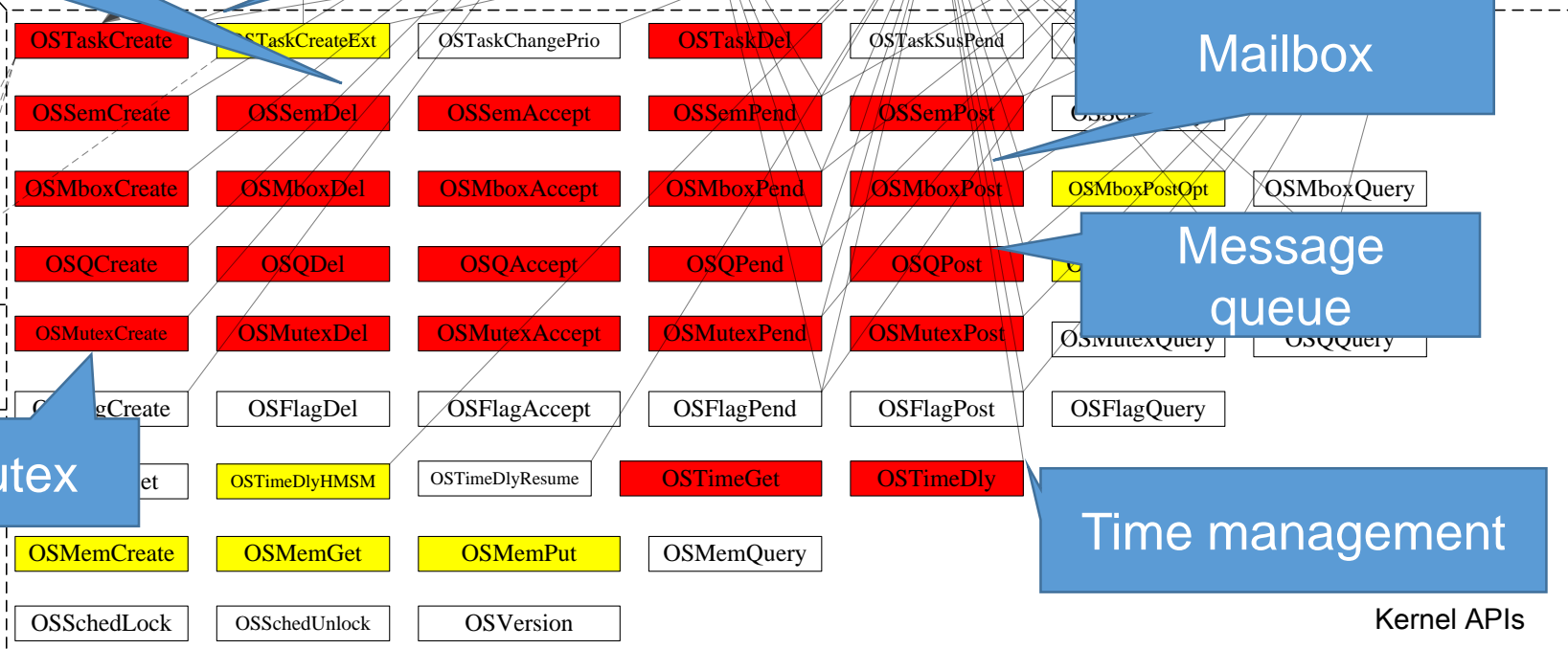
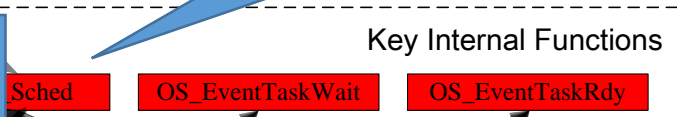
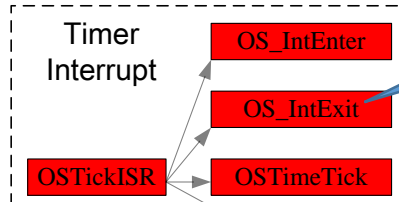
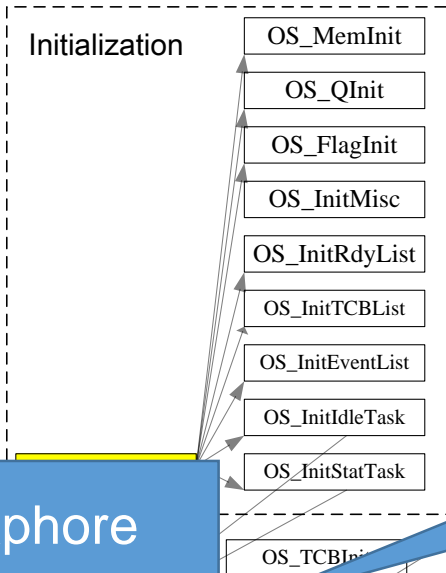
Mailbox

Message queue

Mutex

Time management

Kernel APIs



User Task

OS_TaskIdle

OS_TaskInit

Proving Priority Inversion Freedom

Relational
Assertion

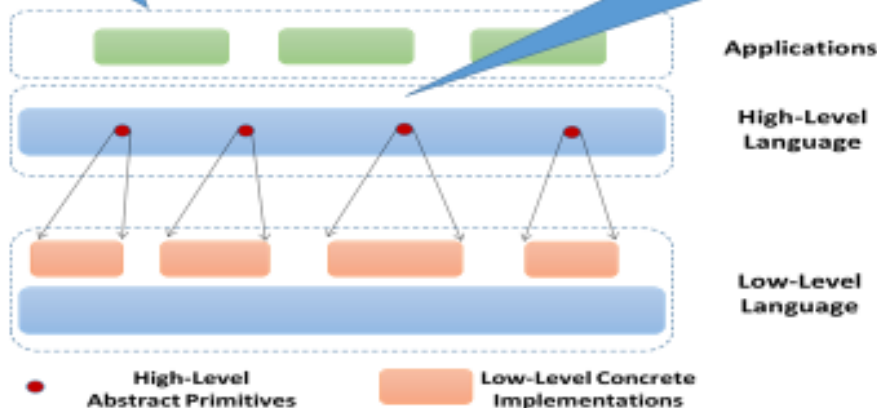
CSL-Style Refinement-Based
Program Logic

Priority inversion
freedom of mutex in
 $\mu\text{C}/\text{OS-II}$

Runtime services and Sys. Props

Whole system properties:
specified as trace properties of
all apps (with high-level views)

Proved at high-level only,
propagated to low-level
through
contextual refinement!



PIF of Mutex

Assembly
Primitives

Operational Semantics
t Switch and Interrupts
-Level Language

framework

Bugs found in μ C/OS-II

- Priority Inversion Freedom in Mutex
 - Use a simplified priority ceiling protocol

Limitation of Mutex



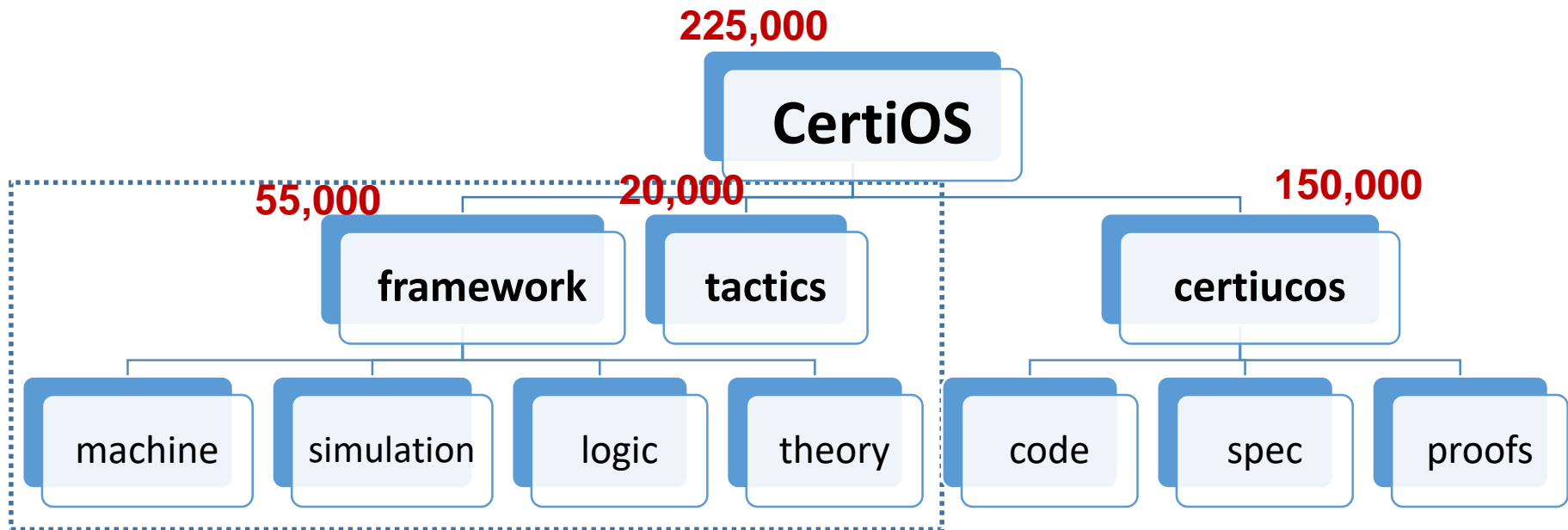
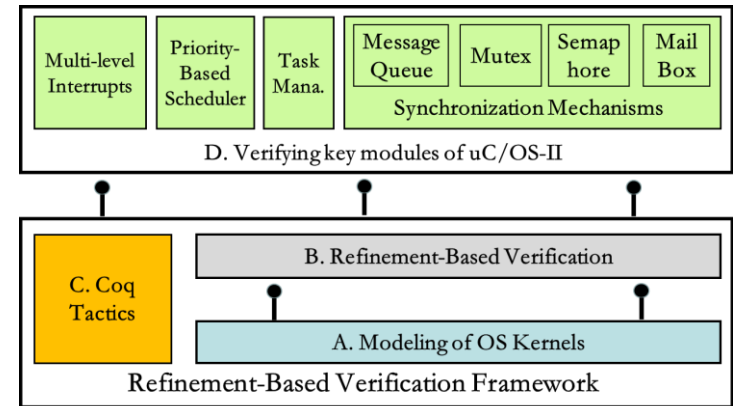
Mutual exclusion semaphores with built-in priority ceiling protocol to prevent priority inversions

- Delivered with complete, clean, consistent, 100% ANSI C source code with in-depth documentation.
- Mutual exclusion semaphores with built-in priority ceiling protocol to prevent priority inversions
- Timeouts on 'pend' calls to prevent deadlocks
- Up to 254 application tasks (1 task per priority level), and unlimited number of kernel objects
- Highly scalable (6K to 24K bytes code space, 1K+ bytes data space)
- Very low interrupt disable time
- Third party certifiable

Bugs found in μ C/OS-II

- Priority Inversion Freedom in Mutex
 - Use a simplified priority ceiling protocol
 - **May cause priority inversion with nested use of mutex!**
 - Fixed in μ C/OS-III
- Concurrency bug (atomicity violation)
INT8U OSTaskChangePrio (INT8U oldprio, INT8U newprio)
 - **May lead to access of invalid pointers**
 - Found in μ C/OS-II v2.52 (the version we verified)
 - Fixed in μ C/OS-II v2.9

Coq Implementations



Time cost: 6 person-years

Components	Cost (py)
Basic framework design and impl.	4
First module: message queue (750 lines of C)	1
the other modules (3000 loc)	1
Total	6

Debugging and fixing framework, specifications, tactics, etc.

Verification can be much faster with stable framework, tools and libraries

<http://staff.ustc.edu.cn/~fuming/research/certiucos/>

Conclusion

- Contextual refinements:
a natural correctness formulation for OS kernels
- Verification framework for preemptive kernels
 - CSL-R: Concurrency refinement + hardware interrupts
- Verification of $\mu\text{C}/\text{OS-II}$
 - Commercial system independently developed by third-party

Limitations & Future Work

- No termination proofs
 - Relatively simple, can be done in logic or using tools
- Assembly and compiler are not verified
 - Ongoing work
- No separate addr. space and isolation
- No real-time properties
- More whole-system properties, in addition to PIF
- Improvements for automation (better tools and libs)

Acknowledgments: Group Members

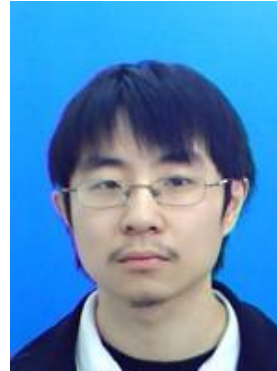
Ming Fu



Fengwei Xu



Xiaoran Zhang



Zhaohui Li



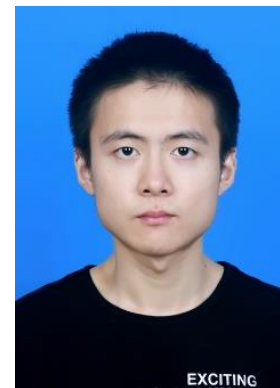
Hui Zhang



Ding Ma



Haibo Gu



Alumni: Jingyuan Cao, Jiebo Ma

Thank you!