

A Practical Verification Framework for Preemptive OS Kernels ^{*} (Technical Report)

Fengwei Xu, Ming Fu^{**}, Xinyu Feng,
Xiaoran Zhang, Hui Zhang, and Zhaohui Li

University of Science and Technology of China

Abstract. We propose a practical verification framework for preemptive OS kernels. The framework models the correctness of API implementations in OS kernels as contextual refinement of their abstract specifications. It provides a specification language for defining the high-level abstract model of OS kernels, a program logic for refinement verification of concurrent kernel code with multi-level hardware interrupts, and automated tactics for developing mechanized proofs. The whole framework is developed for a practical subset of the C language. We have successfully applied it to verify key modules of a commercial preemptive OS $\mu\text{C}/\text{OS-II}$ [2], including the scheduler, interrupt handlers, message queues, and mutexes *etc.* We also verify the priority-inversion-freedom (PIF) in $\mu\text{C}/\text{OS-II}$. All the proofs are mechanized in Coq. To our knowledge, our work is the first to verify the functional correctness of a practical *preemptive* OS kernel with machine-checkable proofs.

1 Introduction

Verifying OS kernels has long been recognized as an important but also extremely challenging task. There have been exciting efforts for OS kernel verification [4, 16, 29, 13] in recent years, but most of them have no or limited support of kernel-level preemption, which allows tasks to be preempted even in kernel mode. This limitation restricts their applicability to real-time systems, where preemptive multitasking is indispensable to achieve real-time guarantees.

Preemptive kernels require explicit invocation of schedulers inside interrupt handlers and careful interrupt management in the kernel code, which make the kernel highly concurrent and complex. In this paper we propose a verification framework for preemptive OS kernels, and show its application in verifying key modules of $\mu\text{C}/\text{OS-II}$ [2], a commercial preemptive real-time multitasking kernel for microprocessors and microcontrollers. The verification is fully mechanized in Coq [1]. To our knowledge, it is the first verification of (key modules of) a *preemptive* OS kernel with machine-checkable proofs. The key contribution of the

^{*} This work is supported in part by grants from National Natural Science Foundation of China (NSFC) under Grant Nos. 61103023, 61229201, 61379039 and 91318301.

^{**} corresponding author (fuming@ustc.edu.cn)

work is to adapt existing theories on interrupt verification [11] and contextual refinement of concurrent programs [18, 20, 27, 26], and integrate them into a framework for real-world preemptive OS kernel verification. Specifically, our work makes the following new contributions:

First, we formulate and verify the correctness of the APIs of OS kernels as *contextual refinement* between their implementations and specifications. Although refinement approaches have been applied in earlier work on OS kernel verification [4, 16, 13], we believe our work is the first to explicitly specify and prove contextual refinement for APIs of a preemptive OS kernel, following recent progress on refinement verification of *concurrent* programs [18, 20, 27, 26]. As we explain in Sec. 2.2, contextual refinement not only serves as a very strong notion of functional correctness of system APIs, but also allows us to prove properties based on the more abstract API specifications and then carry it down to the level of concrete implementations, which makes the verification much simpler than doing proofs directly at the concrete level.

Second, we provide a simple modeling language for specifying kernel primitives. The language strives for balance between abstraction and expressiveness for scheduling. On the one hand, we want the specification to abstract away implementation details. On the other hand, it should provide enough details so that many important properties can be specified at the abstract specification level. Our modeling language provides an abstract **sched** command, allowing us to specify explicitly when the scheduler is invoked in synchronization primitives or interrupt handlers. Semantics of **sched** is parameterized over abstract scheduling policies (*e.g.*, priority-based or round-robin). Expressiveness about these details are necessary to specify system-wide scheduling properties.

Third, we propose a program logic for refinement verification of concurrent kernel programs. The logic supports multi-level nested hardware interrupts and configurable schedulers. It extends concurrent separation logic [22] (CSL) with relational assertions that relate program states at the implementation and the specification levels, as in Liang *et al.* [18, 20]. It also assigns ownership-transfer semantics to interrupt management operations and verify multi-level hardware interrupts in a realistic setting. Different from traditional Hoare-style program logics, whose soundness ensures the semantic interpretation of Hoare-triples, our logic explicitly establishes contextual refinement, which is more useful for establishing abstractions for system APIs, as explained above.

Fourth, our framework is developed for a practical subset of C. It has been successfully applied to verify key APIs of $\mu\text{C}/\text{OS-II}$ [2], including the timer interrupt handler (and a pseudo interrupt handler to demonstrate the support of multi-level interrupts), the scheduler, the time management, and four synchronization mechanisms: message queues, mail boxes, semaphores, and mutexes. It is worth noting that, unlike existing works [4, 16, 29, 13] that are focused on kernels newly developed with verification in mind, we take a *commercial system developed by an independent third-party* and verify the code with minimum modification, which demonstrates the generality and applicability of our framework.

Fifth, we also specify and verify priority inversion freedom (PIF) of $\mu\text{C}/\text{OS-II}$. PIF is a crucial property for real-time systems and is worth verifying in its own right. Moreover, since the specification and verification are done at the level of the abstract model (*i.e.*, specifications) of the kernel, they also help validate our model of system APIs. As we explain above, many important properties cannot be specified if the model is too weak or overly abstract.

Coq proofs are available at

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

2 Background and Overview of Our Work

2.1 Preemptive OS Kernels and Interrupts

In a preemptive OS, execution of a task can be interrupted at any program point (unless interrupts have been disabled) and the control flow can be switched to a different task. To allow preemption, we need these two conditions: (1) enabling interrupts at the kernel level; and (2) invoking the scheduler and context switching inside the interrupt handler.

As shown in Fig. 1, execution of Task A is interrupted and the control is switched to the interrupt handler (step (1)). Instead of returning to Task A directly when the interrupt request has been handled, we may execute the context switch routine (step (4)), which switches the control to another Task B. If we abstract away the interrupt handler, we say the execution of Task A is preempted by Task B in this case. Note that if we disallow one of the above conditions, the only way to switch the control from Task A to task B is to let Task A volunteer to execute `switch` in its code (such as step (11)), resulting in a non-preemptive concurrency model.

Interaction between tasks and hardware interrupts. As we can see, interrupt handling and management are indispensable in preemptive OS kernels. Below we also give a simplified overview of the interrupt mechanism in x86 systems (based on the Intel 8259A interrupt controller).

The CPU has a flag bit `IF` indicating whether interrupts are enabled or not. The `cld/sti` instruction clears/sets the bit to disable/enable interrupts. In 8259A there is a register `isr`, each bit of which corresponds to a hardware interrupt and records if the interrupt is being served or not. Different priority levels are assigned to different sources of interrupts, with level-0 being the highest. When an interrupt request comes, we check `IF` and `isr`. If the interrupts are enabled and there is currently no interrupt of higher or the same priority being served, the request will be served. The corresponding bit in `isr` is set to 1 and the control jumps to the corresponding interrupt handler, shown as step (1) in Fig. 1.

On the invocation of the interrupt handler, the CPU flags (including the `IF` bit) are saved on the stack, and interrupts are disabled automatically. If the programmer enables interrupts again inside the handler, the interrupt handler could be further interrupted by interrupt requests with higher priorities (see

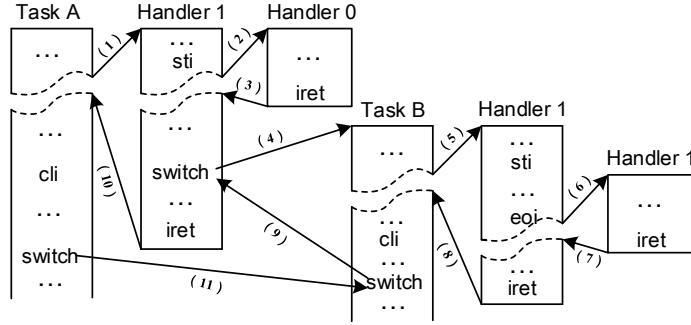


Fig. 1. Tasks and Multi-level Interrupts

step (2), where the level-0 interrupt has a higher priority than level-1 interrupt), causing nested interrupts.

The interrupt handler returns to the program being interrupted using **iret** (step (3)), which also recovers the flags (including the IF bit). Before the handler returns, we need to execute **eo**. The command sends an “end of interrupt” signal to the interrupt controller, which clears the corresponding bit in **isr**. After **eo**, if the interrupt is enabled (**IF** = 1), the interrupt handler could be further interrupted by interrupts at a lower or the same level (see step (6)).

Overview of $\mu C/OS-II$. $\mu C/OS-II$ is a commercial preemptive real-time multi-tasking OS kernel developed by Micrium [2]. The kernel has 6000+ lines of C code and 300+ lines of assembly. It allows a fixed number of tasks, multi-level interrupts, and preemptive priority-based scheduling. The system APIs include “semaphores; event flags; mutual-exclusion semaphores that eliminate unbounded priority inversions; mailboxes; message queues; task, time and timer management; and fixed sized memory block management” [2]. $\mu C/OS-II$ is developed for microprocessors and microcontrollers, and it does not support virtual memory. It has been deployed in many real-world safety critical applications, including avionics (*e.g.*, the Mars Curiosity Rover) and medical equipments.

2.2 Overview of the Verification Framework

An OS kernel hides details of the underlying hardware and provides an abstract programming model for application-level programmers. The implementation of the kernel must ensure that behaviors of user applications in the real machine are consistent with their behaviors under the abstract model [14]. Thus the OS verification can be reduced to verifying refinement between the concrete and abstract programming models.

Contextual refinement as correctness. We consider three entities, the application A , the abstract specifications of the system APIs and interrupt handlers \mathbb{O} , and their concrete implementations O . When system calls are made or interrupts

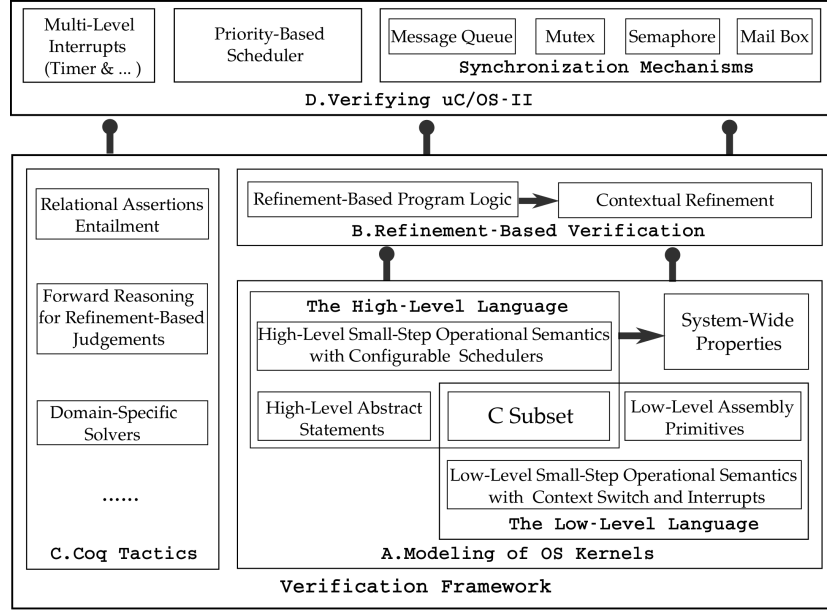


Fig. 2. Structure of the Verification Framework and $\mu\text{C}/\text{OS-II}$ Verification

are handled, routines in O are invoked in the real execution, while in the programmers' mind those in \mathbb{O} are invoked instead at the abstract level. Then the correctness of OS kernels requires O refines \mathbb{O} under *all contexts* A :

$$\forall A. \llbracket A[O] \rrbracket \subseteq \llbracket A[\mathbb{O}] \rrbracket$$

where $\llbracket \cdot \rrbracket$ maps a program P to the set of its observable behaviors. It says that, for all applications, executing the concrete code O does not have more observable behaviors than executing the abstract version \mathbb{O} . In this paper, observable behaviors are defined as finite prefixes of execution traces consisting of observable events, following Liang *et al.* [18].

Contextual refinement is a very strong notion of functional correctness of system APIs since it quantifies over *all* applications. Moreover, it makes verification of system-wide properties simpler. For instance, if we want to verify certain property Φ about a whole system $A[O]$, *i.e.*, Φ holds over every trace in $\llbracket A[O] \rrbracket$, we could prove that it holds over every trace in the superset $\llbracket A[\mathbb{O}] \rrbracket$ instead. Proofs at the abstract level could be much simpler than the concrete level.

The whole verification framework. Figure 2 shows the structure of our verification framework. To model OS kernels and applications, we introduce two languages (in block A), the low-level language for the concrete code implementation and the high-level language for the abstract specification. Above them we have a program logic (in block B) that allows us to prove the low-level kernel implementation contextually refines the high-level specifications. The framework also provides a set of Coq tactics (in block C) to automatically generate and prove

(Addr)	$a \in \text{Int32}$	(Ident)	$id \in \mathbb{Z}$
(FName)	$f \in \mathbb{Z}$	(Integer)	$k \in \text{Int32}$
(Type)	$\tau ::= \text{Tvoid} \mid \text{Tint32} \mid \text{Tptr}(\tau) \mid \text{Tarray}(\tau, n) \mid \text{Tstruct}(id, \mathcal{D}) \mid \dots$		
(OpVal)	$\hat{v} ::= \perp \mid v$	(Value)	$v ::= \text{Vundef} \mid \text{Vnull} \mid \text{Vint}(k) \mid \text{Vptr}(a)$
(ValList)	$\bar{v} ::= \text{nil} \mid v :: \bar{v}$	(ExprList)	$\bar{e} ::= \text{nil} \mid e :: \bar{e}$
(TypeList)	$\mathcal{T} ::= \text{nil} \mid \tau :: \mathcal{T}$	(DeclList)	$\mathcal{D} ::= \text{nil} \mid (id, \tau) :: \mathcal{D}$
(UOP)	$\text{uop} ::= \sim \mid ! \mid \dots$	(BOP)	$\text{bop} ::= + \mid - \mid \gg \mid \ll \mid \& \mid \dots$
(CExpr)	$e ::= x \mid k \mid *e \mid \&e \mid e.id \mid e[e] \mid (\tau)e \mid \text{uop } e \mid e \text{ bop } e$		
(CltStmts)	$d ::= e = e \mid f(\bar{e}) \mid d; d \mid \text{if } (e) \ d \text{ else } d \mid \text{while } (e) \ d \mid \text{return } e \mid \text{print } e \mid \dots$		
(CltFDef)	$cfd ::= (\tau, \mathcal{D}_1, \mathcal{D}_2, d)$	(CltCode)	$A ::= \{f_1 \rightsquigarrow cfd_1, \dots, f_n \rightsquigarrow cfd_n\}$

Fig. 3. The Language for Applications

verification conditions. The $\mu\text{C}/\text{OS-II}$ modules certified in this framework are shown in block D. Below we give details of some of the building blocks.

3 Modeling of the Kernel

As explained above, the correctness of OS kernels is formalized based on three entities — user applications A , the concrete implementation O , and the abstract specification \mathbb{O} . In this section we introduce the programming (model) languages for the three entities.

3.1 The Low-Level Language

The low-level language consists of two parts for implementations of user applications and OS kernels, respectively.

Application language. The application language is shown in Fig. 3. It is a subset of the C language consisting of function calls, pointer operations (except pointer arithmetics), arrays, structs, bit operations, *etc.* The application code A maps function names to their function bodies. The function definition cfd for client consists of the type of the return value, the declaration of parameters, the declaration of local variables and the statements of the function body.

The command $f(\bar{e})$ calls the function f , which could be either an application function in A or an OS API (in O at the low-level or in \mathbb{O} at the high-level, as we explain below).

We use k for 32-bit integers and a for memory addresses (pointers). A value v is either undefined, null, a 32-bit word value or a pointer.

The language supports rich C data types, including the type for the void type (Tvoid), the type of “int” (Tint32), the type of “pointers” ($\text{Tptr}(\tau)$), the type of “array” ($\text{Tarray}(\tau, k)$) and the type of “struct” ($\text{Tstruct}(id, \mathcal{D})$), *etc.*

(LPrim)	$\iota ::= \mathbf{switch} \ x \mathbf{encrt} \mathbf{excrt} \mathbf{eoi} \ k \mathbf{iext} \dots$	
(LStmts)	$s ::= d \iota s ; s \mathbf{while} \ (e) \ s \mathbf{if} \ (e) \ s \ \mathbf{else} \ s$	
(ItrpCode)	$\theta ::= [s_0, \dots, s_{N-1}]$	(LFunDef) $ofd ::= (\tau, \mathcal{D}_1, \mathcal{D}_2, s)$
(ProgUnit)	$\eta ::= \{f_1 \rightsquigarrow ofd_1, \dots, f_n \rightsquigarrow ofd_n\}$	
(LOSCode)	$O ::= (\eta_a, \eta_i, \theta)$	(LProg) $P ::= (A, O)$

Fig. 4. The Languages for Kernel Impl.

(Memory)	$M \in Addr \rightarrow Vaule$	(SymTable)	$G, E \in Var \rightarrow Addr$
(CState)	$\Delta ::= (G, \Pi, M)$	(SymTblSet)	$\Pi ::= \{t_1 \rightsquigarrow E_1, \dots, t_n \rightsquigarrow E_n\}$
(Cont)	$K ::= (\kappa_e, \kappa_s)$	(ExprCont)	$\kappa_e ::= \circ \dots$
(StmtCont)	$\kappa_s ::= \bullet s \cdot \kappa_s \mathbb{S} \cdot \kappa_s (c, \kappa_e, E) \cdot \kappa_s (f, \bar{v}, \bar{e}) \cdot \kappa_s (f \ s \ E) \cdot \kappa_s \dots$		
(CurEval)	$c ::= e s \mathbb{S} \mathbf{fexec}(f, \bar{v}) \mathbf{alloc}(\bar{v}, \mathcal{D}) \mathbf{skip} v \dots$		
(TaskId)	$t \in Addr$	(TaskCode)	$C ::= (c, K)$
(CMem)	$m ::= (G, E, M)$	(TaskPool)	$T ::= \{t_1 \rightsquigarrow C_1, \dots, t_n \rightsquigarrow C_n\}$

Fig. 5. The Common States

An expression e follows the syntax of the C programming language. It is either a constant integer, a program variable, memory reference, deference or standard arithmetic or logical operations over expressions.

A statement d is either an assignment statement, function call statements, a sequence of statements, a branch statement, a loop statement, a return statement or an output statement **print**. The language is reasonably practical because it has been used to implement an executable operating system kernel $\mu\text{C}/\text{OS-II}$. The output command **print** e generates observable event, which is used to define observable event traces needed in our definition of refinement.

Note that the correctness of OS kernels are independent of the implementation language of A . Here we pick the C language for A to simplify the formalization because the applications and the kernel are now implemented in the same language and we do not have to consider the interaction between different languages when defining the whole system ($A[O]$) behaviors.

Low-level language for OS kernels. Figure 4 shows the low-level language for the concrete implementation of OS kernels. Usually the kernels are implemented in C with inline assembly. However, giving semantics directly to C with inline assembly requires us to expose stacks and registers, which makes the semantics overly complex. To avoid this problem, we extend the C statements with assembly primitives ι to encapsulate the assembly code. Semantics of these primitives will be given below.

switch x switches to the target task x . **encrt** enters a critical region by disabling interrupts. It also saves the old IF onto the stack to allow nested critical regions. Note we use ie to model the IF flag and abstract away other bits in the hardware EFLAGS register. **excrt** exits the current critical region by

popping the stack to recover ie . Since we hide stacks in our state model, we use an abstract stack cs to save the historical ie bits (see Fig. 6, which is explained below). **eo** i k clears the k -th bit in isr , indicating that the k -th interrupt is no longer in service. **ie****xt** enables interrupts and returns to the interrupted program.

The kernel implementation O consists of the system API implementation η_a , the internal functions η_i and the interrupt handlers θ . The internal functions are called only by code in η_a or θ . θ is a sequence of N interrupt handlers, where N is the maximum number of interrupts we support. The handler with the lower identifier has the higher priority. Then a complete low-level program P is defined as a pair of the application code A and the kernel code O .

Common machine states. As shown in Fig. 5, we present the common machine states for the two levels. The memory M is modeled as a partial function from addresses to values. The global symbol table G and the local symbol table E map program variables to addresses. Note that we use a flat memory model to simplify the presentation. The basic memory operations follow the block-based memory model in CompCert [17]. Π maps the task identifiers to their local symbol tables. Δ consists of the global symbol table G , the set of local symbol tables Π and the memory M .

We give small-step operational semantics at the two levels. For each step, the processor picks the continuation of the current task and executes its current command or expression. To model fine-grained concurrency, both commands and *expressions* could be executed in multiple steps, where each step corresponds to the granularity of a single machine instruction (as in CompCertTSO [24], but we use the sequential consistent model instead of the x86-TSO memory model).

The expression and statement continuations are presented in Fig. 5. For the expression continuations κ_e , they can be \circ which means that there is nothing left to be evaluated, or some other standard cases.

For the statement continuations κ_s , \bullet means that there is nothing left to be done. When the currently running task is interrupted, we use $(c, \kappa_e, E) \cdot \kappa_s$ to save the execution context for the current task. Since the interrupt may happen when evaluating an expression, we need to record the current evaluation c , the expression continuation κ_e and the current local symbol table E for resuming the execution context in the future. We use $(f, \bar{v}, \bar{e}) \cdot \kappa_s$ to save intermediate results of calculating function arguments. $(f, s, E) \cdot \kappa_s$ is used to save the context when doing function calls.

We also introduce some runtime statements for defining the small-step operational semantics. For instance, **fexec** (f, \bar{v}) is an intermediate statement for calling a function. **alloc** (\bar{v}, \mathcal{D}) is used to do memory allocations for local variables and parameters in a function. More details about the usage of these statements can be seen in Fig 8, where we give some key rules of operational semantics for the low-level language.

Low-level machine states. The language is concurrent, with multiple continuations (*i.e.*, control stacks) in the state, each corresponding to a task. All tasks share memory, but each has its own local variables and local interrupt states

$(BitVal)$	$b, ie \in \{0, 1\}$	$(ISRReg)$	$isr ::= [b_0, \dots, b_{N-1}]$
$(CrtStk)$	$cs ::= nil \mid ie :: cs$	$(ItrpStk)$	$is ::= nil \mid k :: is$
$(ItrpTaskSt)$	$\delta ::= (ie, is, cs)$	$(ItrpSt)$	$\pi ::= \{t_1 \rightsquigarrow \delta_1, \dots, t_n \rightsquigarrow \delta_n\}$
$(LOsFullSt)$	$\Lambda ::= (\Delta, isr, \pi)$	$(TaskLocalSt)$	$\sigma ::= (m, isr, \delta)$
$(LWorld)$	$W ::= (P, T, \Delta, \Lambda, t)$		

Fig. 6. The Low-level Machine States

(see δ in Fig. 6, which is explained below). We also separate the program state (including memory and variables) into two disjoint parts, one for the application code A and the other for the kernel code O . The only way for A to access kernel states is to call system APIs in O , and O cannot access application states.

As explained in Sec. 2.1, ie is a boolean flag used to turn on/off interrupts. isr is a sequence of boolean flags, one for each interrupt. The stack cs records the historical values of ie , which are pushed whenever the execution enters a critical region. It is used to support nested critical regions. The task-local stack is records the sequence of interrupts that interrupt the execution of *this task*. It is auxiliary data introduced for verification purpose. Then the task-local interrupt status δ is defined as a triple (ie, is, cs) . π records the δ of each task. The kernel-level state Λ consists of the general C state Δ , the global isr register and the set π of task-local interrupt status.

The whole program configuration W now consists of the task pool T , the client state Δ , the kernel state Λ , and the identifier t of the current task. Note that W contains two pieces of Δ , one for user applications (clients), the other inside Λ for the kernel. Separating the data into two parts prevents user applications from accessing kernel data. Applications trying to access data unavailable in the client Δ will trigger a runtime error in our operational semantics.

Low-level operational semantics. We give the low-level operational semantics in Fig. 8. The low-level program steps denoted as “ $P \vdash W \xRightarrow{L} W'$ ”, may execute a regular command in a task (the PTASK rule), or execute **switch** x to do context switch (the PSW rule), or be interrupted and transfer the control to the corresponding interrupt handler (the PITRP rule). When executing a regular command in a task, it may either belong to the kernel (the TKERNEL rule) or the client (the TCLT rule).

The PTASK rule is used to lift task-local steps of the current task to program steps. We use $\Lambda|_t$ to project the task-local data σ of t from Λ (σ defined in Fig. 5), then the program configuration is updated according to the execution of the task-local step. Here $\text{UPDTS}(\Lambda, t, \sigma')$ (defined in Fig.7) updates the local data of t in Λ with the new σ' .

The assembly implementation of the context switch routine is abstracted into the primitive **switch** x . It switches the execution from the current task to the target task x , where x stores the task identifier. The PSW rule simply resets the current thread identifier and updates the global variable `OSTCBCur` accordingly.

$$\begin{aligned}
\llbracket \kappa_s \rrbracket_c &\stackrel{\text{def}}{=} \begin{cases} \perp & \text{if } \kappa_s = \bullet \\ \kappa_s & \text{if } \kappa_s = (f, s, E) \cdot \kappa'_s \\ \llbracket \kappa'_s \rrbracket_c & \text{otherwise} \end{cases} & \llbracket \kappa_s \rrbracket &\stackrel{\text{def}}{=} \begin{cases} \perp & \text{if } \kappa_s = \bullet \\ \kappa_s & \text{if } \kappa_s = (c, \kappa_e, E) \cdot \kappa'_s \\ \llbracket \kappa'_s \rrbracket & \text{otherwise} \end{cases} \\
f \perp g &\stackrel{\text{def}}{=} \text{dom}(f) \cap \text{dom}(g) = \emptyset & f \uplus g &\stackrel{\text{def}}{=} \begin{cases} f \cup g & \text{iff } f \perp g \\ \text{undef} & \text{otherwise} \end{cases} \\
\text{InOS}(C, (A, O)) &\stackrel{\text{def}}{=} \exists c, \kappa_s. C = (c, (_, \kappa_s)) \wedge ((\exists f. f \in \text{dom}(O.\eta_a \uplus O.\eta_i) \wedge \\
&\quad (c = \mathbf{fexec}(f, _) \vee \llbracket \kappa_s \rrbracket_c = (f, _, _) \cdot _)) \vee \llbracket \kappa_s \rrbracket \neq \perp) \\
(G, \Pi, M)|_t &= (G, E, M) \stackrel{\text{def}}{=} \Pi(t) = E \\
(\Delta, isr, \pi)|_t &= (m, isr, \delta) \stackrel{\text{def}}{=} \Delta|_t = m \wedge \pi(t) = \delta \\
\Lambda' = \text{UPDTS}(\Lambda, t, \sigma') &\stackrel{\text{def}}{=} \Lambda'|_t = \sigma' \wedge \forall t' \neq t. \Lambda'|_{t'} = \Lambda|_{t'} \\
\Delta' = \text{UPDCS}(\Delta, t, m') &\stackrel{\text{def}}{=} \Delta'|_t = m' \wedge \forall t' \neq t. \Delta'|_{t'} = \Delta|_{t'} \\
\llbracket x \rrbracket_{(G, E, M)} = t' &\stackrel{\text{def}}{=} \exists a. (E(x) = a \wedge M(a) = t') \vee \\
&\quad (G(x) = a \wedge x \notin \text{dom}(E) \wedge M(a) = t') \\
\llbracket x \rrbracket_{(t, ((G, \Pi, M), isr, \pi))} = t' &\stackrel{\text{def}}{=} \llbracket x \rrbracket_{(G, \Pi(t), M)} = t' \\
((G, \Pi, M'), isr, \pi) = \text{UPDG}(((G, \Pi, M), isr, \pi), \text{OSTCBCur}, t') &\stackrel{\text{def}}{=} \\
&\quad \exists a. G(\text{OSTCBCur}) = a \wedge M' = M\{a \rightsquigarrow t'\}
\end{aligned}$$

Fig. 7. Auxiliary Definitions

The PITRP rule says that the task t can be interrupted by a level- k interrupt request if ie is 1 (thus we are not in critical regions and cs must be `nil`) and there is currently no interrupt of higher or the same priority being served, according to isr . Then we switch the control to the interrupt handler $\theta(k)$, and saves the execution context (c, κ_e, E) onto the statement continuation (see Fig. 8). We also set the k -th bit of isr , clear the ie bit, and push k onto the is stack.

We use “ $P \vdash (C, \Delta, \sigma) \xrightarrow{L} (C', \Delta', \sigma')$ ” to define task-local semantics of the assembly primitives for interrupt management. The TKERNEL rule means to execute a kernel command. It checks whether it is executing the kernel code using $\text{InOS}(C, P)$ (defined in Fig.7 by checking the current continuation). The TCLT rule executes the client code. Here $\text{UPDCS}(\Delta, t, m')$ (defined in Fig.7) updates the local data of t in Δ with the new m' .

The assembly primitives ι except **switch** are all related to interrupts management and handling. To model their semantics, we introduce interrupt states in the state model, as shown at Fig. 6. The *global* register isr is shared by all tasks. It models the isr register in 8259A interrupt controller, as explained in Sec. 2.1. In addition, there are *local* interrupt states δ for each task. It contains a local copy ie of the IF flag in the EFLAGS register (see Sec. 2.1) recording whether interrupts are enabled, a stack cs consisting of the historical values of ie to support nested critical regions, and another stack is recording the sequence

$$\boxed{P \vdash W \Longrightarrow W'}$$

$$\frac{\Lambda|_t = \sigma \quad P \vdash (C, \sigma, \Delta) \xrightarrow{L} (C', \sigma', \Delta') \quad T(t) = C \quad T' = T\{t \rightsquigarrow C'\} \quad \Lambda' = \text{UPDTS}(\Lambda, t, \sigma')}{P \vdash (T, \Delta, \Lambda, t) \Longrightarrow (T', \Delta', \Lambda', t)} \text{ (PTASK)}$$

$$\frac{T(t) = (\mathbf{switch} \ x, K) \quad \llbracket x \rrbracket_{(t, \Lambda)} = t' \quad \Lambda' = \text{UPDG}(\Lambda, \text{OSTCBCur}, t')}{P \vdash (T, \Delta, \Lambda, t) \Longrightarrow (T\{t \rightsquigarrow (\mathbf{skip}, K)\}, \Delta, \Lambda', t')} \text{ (PSW)}$$

$$\frac{P = (A, (\eta_a, \eta_i, \theta)) \quad T(t) = (c, (\kappa_e, \kappa_s)) \quad \Lambda|_t = ((G, E, M), \text{isr}, (1, \text{is}, \text{nil})) \quad \forall k'. k' \leq k \rightarrow \text{isr}(k') = 0 \quad C' = (\theta(k), (\circ, (c, \kappa_e, E) \cdot \kappa_s)) \quad T' = T\{t \rightsquigarrow C'\} \quad \sigma' = ((G, \emptyset, M), \text{isr}\{k \rightsquigarrow 1\}, (0, k :: \text{is}, \text{nil})) \quad \Lambda' = \text{UPDTS}(\Lambda, t, \sigma')}{P \vdash (T, \Delta, \Lambda, t) \Longrightarrow (T', \Delta, \Lambda', t)} \text{ (PITRP)}$$

$$\boxed{P \vdash (C, \Delta, \sigma) \xrightarrow{L} (C', \Delta', \sigma')}$$

$$\frac{\text{lnOS}(C, P) \quad P \vdash (C, \sigma) \bullet \xrightarrow{L} (C', \sigma')}{P \vdash (C, \Delta, \sigma) \xrightarrow{L} (C', \Delta, \sigma')} \text{ (TKERNEL)}$$

$$\frac{\neg \text{lnOS}(C, (A, (\eta_a, \eta_i, \theta))) \quad \Delta|_t = m \quad A \uplus \eta_a \vdash (C, m) \mapsto (C', m') \quad \Delta' = \text{UPDCS}(\Delta, t, m')}{(A, (\eta_a, \eta_i, \theta)) \vdash (C, \Delta, \Lambda) \xrightarrow{L} (C', \Delta', \Lambda)} \text{ (TCLT)}$$

$$\frac{\sigma = (m, \text{isr}, (ie, \text{is}, cs)) \quad \sigma' = (m, \text{isr}, (0, \text{is}, ie :: cs))}{P \vdash ((\mathbf{encrt}, K), \sigma) \bullet \xrightarrow{L} ((\mathbf{skip}, K), \sigma')} \text{ (ENTERCRT)}$$

$$\frac{\sigma = (m, \text{isr}, (ie, \text{is}, ie' :: cs)) \quad \sigma' = (m, \text{isr}, (ie', \text{is}, cs))}{P \vdash ((\mathbf{excrt}, K), \sigma) \bullet \xrightarrow{L} ((\mathbf{skip}, K), \sigma')} \text{ (EXITCRT)}$$

$$\frac{0 \leq k < N \quad \sigma = (m, \text{isr}, (ie, \text{is}, cs)) \quad \sigma' = (m, \text{isr}\{k \rightsquigarrow 0\}, (ie, \text{is}, cs))}{P \vdash ((\mathbf{eoi} \ k, K), \sigma) \bullet \xrightarrow{L} ((\mathbf{skip}, K), \sigma')} \text{ (EOI)}$$

$$\frac{\sigma = ((G, E, M), \text{isr}, (ie, k :: \text{is}, cs)) \quad \sigma' = ((G, E', M), \text{isr}, (1, \text{is}, cs)) \quad \lfloor \kappa_s \rfloor = (c, \kappa_e, E') \cdot \kappa'_s}{P \vdash ((\mathbf{iext}, (\circ, \kappa_s)), \sigma) \bullet \xrightarrow{L} ((c, (\kappa_e, \kappa'_s)), \sigma')} \text{ (IRET)}$$

$$\frac{\eta_i \vdash (C, m) \mapsto (C', m')}{(A, (\eta_a, \eta_i, \theta)) \vdash (C, (m, \text{isr}, \delta)) \bullet \xrightarrow{L} (C', (m', \text{isr}, \delta))} \text{ (KCSTEP)}$$

$$\boxed{\eta \vdash (C, m) \mapsto (C', m')}$$

$$\frac{}{\eta \vdash ((f(\text{nil}), (\circ, \kappa_s)), m) \mapsto ((\mathbf{fexec}(f, \text{nil}), (\circ, \kappa_s)), m)} \text{ (FN)}$$

$$\frac{}{\eta \vdash ((f(e :: \bar{e}), (\circ, \kappa_s)), m) \mapsto ((e, (\circ, (f, \text{nil}, \bar{e}) \cdot \kappa_s)), m)} \text{ (FA)}$$

$$\frac{}{\eta \vdash ((v, (\circ, (f, \bar{v}, (e :: \bar{e})) \cdot \kappa_s)), m) \mapsto ((e, (\circ, (f, (v :: \bar{v}), \bar{e}) \cdot \kappa_s)), m)} \text{ (FEVAL)}$$

$$\frac{}{\eta \vdash ((v, (\circ, (f, \bar{v}, \text{nil}) \cdot \kappa_s)), m) \mapsto ((\mathbf{fexec}(f, v :: \bar{v}), (\circ, \kappa_s)), m)} \text{ (FENTER)}$$

$$\frac{\eta(f) = (\tau, \mathcal{D}_1, \mathcal{D}_2, s) \quad m = (G, E, M) \quad m' = (G, \emptyset, M)}{\eta \vdash ((\mathbf{fexec}(f, \bar{v}), (\circ, \kappa_s)), m) \mapsto ((\mathbf{alloc}(\bar{v}, \text{rev}(\mathcal{D}_1) + \mathcal{D}_2), (\circ, (f, s, E) \cdot \kappa_s)), m')} \text{ (FALLOC)}$$

$$\frac{}{\eta \vdash ((\mathbf{alloc}(\text{nil}, \text{nil}), (\circ, (f, s, E) \cdot \kappa_s)), m) \mapsto ((s, (\circ, (f, s, E) \cdot \kappa_s)), m)} \text{ (FBODY)}$$

Fig. 8. The Low-level Operational Semantics

of interrupts that interrupt the execution of *the task*. The stack *is* is auxiliary data introduced mainly for verification purposes. π records the δ of each task.

encrt enters a critical region by disabling interrupts (*i.e.*, clearing the *ie* bit using **cli**). It also saves the old *ie* onto the *cs* stack. **excrt** exits the critical region by popping off the top value on *cs* and using it to restore *ie* (executing **sti** if the value is 1).

Interrupt requests may arrive non-deterministically after each step if *ie* = 1. A level-*k* request is served only if there is no request at higher or the same level being served (*i.e.*, $\forall k'. k' \leq k \rightarrow isr(k') = 0$). Then the processor clears *ie*, sets *isr*(*k*) to 1, pushes the number *k* onto the logical stack *is*, saves the execution context and the local variables onto the abstract control stack (*i.e.*, the continuation), and finally jumps to the interrupt handler $\theta(k)$.

eo *k* clears the *k*-th bit in *isr*, indicating that the *k*-th interrupt is no longer in service. **ixt** is an abstraction of the **iret** instruction. It resets the *ie* bit to 1 to enable interrupts, pops out the topmost interrupt number on the *is* stack, and returns to the interrupted program. The **IRET** rule pops the execution context from the statement continuation, and sets the *ie* bit. We use $\lfloor \kappa_s \rfloor$ defined in Fig.7 to pop the execution context from the statements continuations.

In addition, we use “ $\eta \vdash (C, m) \mapsto (C', m')$ ” to define the operational semantics for common C steps. Here we only give the **FN**, **FA**, **FEVAL**, **FENTER**, **FALLOC** and **FBODY** rules for executing a function call, other rules for standard C semantics are in our Coq code and omitted here. When invoking a function, it firsts uses the **FA** rule to evaluate the expressions of function arguments. Then the **FEVAL** rule is applied for evaluating the next expression. After all the arguments are evaluated to a value list, we use **FALLOC** to allocate memory blocks for local variables and arguments. Finally, we apply the **FBODY** rule to execute the function body after finishing local allocations. Note that the dynamical statement **fexec**(*f*, *v*) is the execution boundary between client-steps and kernel-steps when client code invokes a kernel API.

Encoding C code of $\mu C/OS-II$ in Coq. We do deep encoding of the code by defining the abstract syntax tree of the C language (subset) as an inductive datatype in Coq, and manually write the code as an expression of this type. With the help of Coq notations, the syntactic representations of C code in Coq look similar as the original C code. Actually the encoding could be done by an automatic transformer. Figure 9 (b) shows our encoding in Coq for the C source code of the scheduler of $\mu C/OS-II$ (Fig. 9 (a)). All the code definitions are located in the directory of our Coq source code “CertiOS/certiucos/code/”.

3.2 The High-Level Specification Language

Viewing from the aspect of application programmers, we model the OS kernel as an extended C language with multi-tasking and system calls. As explained above, the C language is used to implement user applications *A*, and the system calls invoke an abstract version of system routines in \mathbb{O} , which are implemented using a simple specification language. Correspondingly, the low-level concrete

```

void OS_Sched (void)
{
    INT8U      y;

    OS_ENTER_CRITICAL();
    y          = OSUnMapTbl[OSRdyGrp];
    OSPrioHighRdy = (INT8U)((y << 3) + OSUnMapTbl[OSRdyTbl[y]]);
    if (OSPrioHighRdy != OSPrioCur) {
        OSTCBHighRdy = OSTCBPrioTbl[OSPrioHighRdy];
        OSCtxSwCtr++;
        OSCtxSw();
    }
    OS_EXIT_CRITICAL();
    return;
}

```

(a) The C code of Scheduler

```

Definition OS_Sched_impl :=
Void ·OS_Sched·(L U)··{
    L y @ Int8U;

    ENTER_CRITICAL;;
    y' = OSUnMapTbl'[OSRdyGrp]';
    OSPrioHighRdy' = ((y' << '3) + OSUnMapTbl'[OSRdyTbl'[y]]);
    IF (OSPrioHighRdy' != OSPrioCur') {
        OSTCBHighRdy' = OSTCBPrioTbl'[OSPrioHighRdy'];
        ++ OSCtxSwCtr';
        SWITCH
    };
    EXIT_CRITICAL;;
    RETURN
}..

```

(b) The Coq encoding of Scheduler

Fig. 9. C Source code vs. Coq encoding

representation of kernel states is modeled as algebraic abstract states at the high level. This section presents the high-level language and its semantics.

As shown in Fig. 10, the whole high-level program \mathbb{P} consists of the application code A and the abstract specification of the kernel \mathbb{O} . The application code A is the same as in the low-level language (see Fig. 4). \mathbb{O} contains the specifications φ for kernel APIs, ε for interrupt handlers, and χ for the scheduler. The high-level program configuration \mathbb{W} consists of the task pool T , the client state Δ , and the abstract kernel state Σ .

Programmers at this level have no control over interrupts (*e.g.*, enabling or disabling interrupts). Always enabled, interrupts are modeled implicitly as abstract external events that may occur non-deterministically at any program points. Handlers of the events are also specified as ε in \mathbb{O} . At the high level an incoming level- k event is always handled by executing $\varepsilon(k)$.

$(HStmts)$	$s ::= \mathbf{sched} \mid \gamma(\bar{v}) \mid \mathbf{assert} \ b \mid \mathbf{end} \ \hat{v} \mid s_1; s_2 \mid s_1 + s_2$	
$(HAPISet)$	$\varphi ::= \{f_1 \rightsquigarrow s_1, \dots, f_n \rightsquigarrow s_n\}$	$(HEvtSet) \ \varepsilon ::= [s_0, \dots, s_{N-1}]$
$(HSched)$	$\chi \in HAbsSt \rightarrow TaskId \rightarrow Prop$	$(HBEpr) \ \flat \in HAbsSt \rightarrow Prop$
$(HOSCode)$	$\mathbb{O} ::= (\varphi, \varepsilon, \chi)$	$(HProg) \ \mathbb{P} ::= (A, \mathbb{O})$
$(HWorld)$	$\mathbb{W} ::= (\mathbb{P}, T, \Delta, \Sigma)$	

Fig. 10. High-Level Spec. Language

$(HAbsSt)$	$\Sigma ::= \{a_1 \rightsquigarrow \Omega_1, \dots, a_n \rightsquigarrow \Omega_n\}$	$(HDataNm)$	$a ::= \mathbf{tcbls} \mid \mathbf{ecbls} \mid \mathbf{ctid} \mid \dots$
$(HData)$	$\Omega ::= \alpha \mid \beta \mid t \mid \dots$	$(HEvtId)$	$eid \in Addr$
$(HTCBLs)$	$\alpha ::= \{t_1 \rightsquigarrow (pr_1, ts_1), \dots, t_n \rightsquigarrow (pr_n, ts_n)\}$		
$(HECBLs)$	$\beta ::= \{eid_1 \rightsquigarrow ed_1, \dots, eid_n \rightsquigarrow ed_n\}$		
$(Priority)$	$pr \in int32$	$(WaitType)$	$wt ::= \mathbf{mtx}(eid) \mid \mathbf{tm} \mid \dots$
$(WaitQ)$	$Q \in \mathbf{nil} \mid t :: Q$	$(HStatus)$	$ts ::= \mathbf{rdy} \mid \mathbf{wait}(wt, k)$
$(MtxOwner)$	$w ::= \perp \mid (t, pr)$	$(HECBData)$	$ed ::= \mathbf{mutex}(pr, w) \mid \dots$

Fig. 11. The High-level Abstract Machine

The system APIs and interrupt handlers are specified as an abstract statement s , which forms a simple but expressive specification language. **sched** does scheduling. Its semantics is determined by the abstract scheduler specification χ . As defined in Fig. 10, χ is a binary relation between abstract states and task identifiers. That is, given an abstract state Σ (defined at the bottom of Fig. 10), χ finds a related task identifier as the next task to execute. Note that χ is a relation instead of a function, therefore the abstract scheduler does not have to be deterministic. Since χ is provided as part of the kernel specification, the semantics of **sched** in our language is configurable. Specifying details of the scheduling policies (instead of using a more abstract non-deterministic scheduler that may pick *any* task) allows us to specify and verify scheduling properties such as PIF at the high level.

$\gamma(\bar{v})$ is a meta-level relation (defined in Coq) that takes \bar{v} as arguments and maps an abstract state to another. Users can instantiate it to specify any *atomic* transitions over abstract states. **assert** b asserts that the predicate b holds over the current abstract state. **end** \hat{v} represents the end of abstract APIs with optional return values or interrupt handlers. $s_1; s_2$ and $s_1 + s_2$ are statements for sequential composition and non-deterministic choices respectively.

Abstract states. The kernel state is represented as the abstract state Σ at the high level. As defined at Fig. 11, Σ is a mapping from names a to the abstract data Ω . Here \mathbf{tcbls} is the name for the high-level abstract TCB list α , which maps task identifiers to abstract tasks, including the priority pr (a natural number), the task status (ready, waiting, *etc.*) and so on, depending on the low-level implementations. \mathbf{ctid} is the name for the current task identifier t .

$$\boxed{\mathbb{P} \vdash W \Rightarrow W'}$$

$$\frac{\frac{\Sigma(\text{ctid}) = t \quad T(t) = C \quad T' = T\{t \rightsquigarrow C'\}}{\mathbb{P} \vdash (C, \Delta, \Sigma) \xrightarrow{H} (C', \Delta', \Sigma')} \quad (\text{HTASK})}{\mathbb{P} \vdash (T, \Delta, \Sigma) \Rightarrow (T', \Delta', \Sigma')}$$

$$\frac{\frac{\mathbb{P} = (A, (\varphi, \varepsilon, \chi)) \quad \Sigma(\text{ctid}) = t \quad \chi \Sigma t'}{T(t) = (\text{sched}; \mathfrak{s}, K) \quad T' = T\{t \rightsquigarrow (\mathfrak{s}, K)\}} \quad (\text{SCHD})}{\mathbb{P} \vdash (T, \Delta, \Sigma) \Rightarrow (T', \Delta, \Sigma\{\text{ctid} \rightsquigarrow t'\})}$$

$$\frac{\Sigma(\text{ctid}) = t \quad T(t) = (c, (\kappa_e, \kappa_s)) \quad T' = T\{t \rightsquigarrow (\varepsilon(k), (\circ, (c, \kappa_e, \emptyset) \cdot \kappa_s))\}}{(A, (\varphi, \varepsilon, \chi)) \vdash (T, \Delta, \Sigma) \Rightarrow (T', \Delta, \Sigma')} \quad (\text{PEVENT})$$

$$\boxed{\mathbb{P} \vdash (C, \Delta, \Sigma) \xrightarrow{H} (C', \Delta, \Sigma')}$$

$$\frac{(\mathfrak{s}, \Sigma) \bullet \xrightarrow{H} (\mathfrak{s}', \Sigma')}{\mathbb{P} \vdash ((\mathfrak{s}, K), \Delta, \Sigma) \xrightarrow{H} ((\mathfrak{s}, K), \Delta, \Sigma')} \quad (\text{INAPI})$$

$$\frac{\frac{\Delta|_t = m \quad A \vdash (C, m) \mapsto (C', m')}{T' = T\{t \rightsquigarrow C'\} \quad \Delta' = \text{UPDCS}(\Delta, t, m')} \quad (\text{HTCLT})}{(A, \emptyset) \vdash (C, \Delta, \Sigma) \xrightarrow{H} (C', \Delta', \Sigma')}$$

$$\frac{\varphi(f) = \omega \quad C = (\text{fexec}(f, \bar{v}), K)}{(A, (\varphi, \varepsilon, \chi)) \vdash (C, \Delta, \Sigma) \xrightarrow{H} ((\omega \bar{v}, K), \Delta, \Sigma)} \quad (\text{ENAPI})$$

$$\frac{}{\mathbb{P} \vdash ((\text{end } v, K), \Delta, \Sigma) \xrightarrow{H} ((v, K), \Delta, \Sigma)} \quad (\text{ENDAPI1})$$

$$\frac{}{\mathbb{P} \vdash ((\text{end } , K), \Delta, \Sigma) \xrightarrow{H} ((\text{skip}, K), \Delta, \Sigma)} \quad (\text{ENDAPI2})$$

$$\frac{C = (\text{end } , (\circ, (c, \kappa_e, \emptyset) \cdot \kappa_s))}{\mathbb{P} \vdash (C, \Delta, \Sigma) \xrightarrow{H} ((c, (\kappa_e, \kappa_s)), \Delta, \Sigma)} \quad (\text{ENDEVT})$$

$$\boxed{(\mathfrak{s}, \Sigma) \bullet \xrightarrow{H} (\mathfrak{s}', \Sigma')}$$

$$\frac{\frac{\gamma \bar{v} \Sigma (\hat{v}, \Sigma') \quad \Sigma \perp \Sigma_f \quad \text{dom}(\Sigma(\text{tbls})) = \text{dom}(\Sigma'(\text{tbls}))}{\text{dom}(\Sigma) = \text{dom}(\Sigma') \quad \Sigma(\text{ctid}) = \Sigma'(\text{ctid})}}{(\gamma(\bar{v}), \Sigma \cup \Sigma_f) \bullet \xrightarrow{H} (\text{end } \hat{v}, \Sigma' \cup \Sigma_f)} \quad (\text{PRIM})$$

$$\frac{\mathbb{b} \Sigma \quad \Sigma \perp \Sigma_f}{(\text{assert } \mathbb{b}, \Sigma \cup \Sigma_f) \bullet \xrightarrow{H} (\text{end } \perp, \Sigma \cup \Sigma_f)} \quad (\text{ASSERT})$$

$$\frac{}{(\text{end } \hat{v}; \mathfrak{s}, \Sigma) \bullet \xrightarrow{H} (\mathfrak{s}, \Sigma)} \quad \frac{(\mathfrak{s}_1, \Sigma) \bullet \xrightarrow{H} (\mathfrak{s}'_1, \Sigma')}{(\mathfrak{s}_1; \mathfrak{s}_2, \Sigma) \bullet \xrightarrow{H} (\mathfrak{s}'_1; \mathfrak{s}_2, \Sigma')}$$

$$\frac{}{(\mathfrak{s}_1 + \mathfrak{s}_2, \Sigma) \bullet \xrightarrow{H} (\mathfrak{s}_1, \Sigma)} \quad \frac{}{(\mathfrak{s}_1 + \mathfrak{s}_2, \Sigma) \bullet \xrightarrow{H} (\mathfrak{s}_2, \Sigma')}$$

Fig. 12. The High-level Operational Semantics

```

void OSTimeDly (Int16u ticks) :
1  if (ticks > 0){
2    OS_ENTER_CRITICAL();
3    if((OSTCBCur->OSTCBStat==OS_STAT_RDY)&&
        (OSTCBCur->OSTCBDly==0)){
4      OSRdyTbl[OSTCBCur->OSTCBY]=
        OSRdyTbl[OSTCBCur->OSTCBY] & (OSTCBCur->OSTCBBitX);
5      if(OSRdyTbl[OSTCBCur->OSTCBY]==0){
6        OSRdyGrp= OSRdyGrp & (OSTCBCur->OSTCBBitY)}
7      OSTCBCur->OSTCBDly=ticks;
8      OS_EXIT_CRITICAL();
9      OSSched();
10     }else{ OS_EXIT_CRITICAL();} }
11  return;

OSTimeDly:  $\mathbb{s}_{\text{dly}} \stackrel{\text{def}}{=} (\gamma_{\text{err}}(\text{ticks}) + (\gamma_{\text{dly}}(\text{ticks}); \text{sched}))$ 
 $\gamma_{\text{err}}(\text{ticks}) \stackrel{\text{def}}{=} \lambda \Sigma, (\hat{v}, \Sigma'). \Sigma = \Sigma' \wedge \text{ticks} = 0$ 
 $\gamma_{\text{dly}}(\text{ticks}) \stackrel{\text{def}}{=} \lambda \Sigma, (\hat{v}, \Sigma'). \text{ticks} > 0 \wedge (\exists t, pr. \Sigma(\text{ctid}) = t \wedge \Sigma(\text{tcbls})(t) = (pr, \text{rdy}) \wedge$ 
 $\Sigma' = \Sigma\{\text{tcbls} \rightsquigarrow \{t \rightsquigarrow (pr, \text{wait}(wt, \text{ticks}))\}\})$ 

```

Fig. 13. Specification Code for OSTimeDly

wt represents the type of waiting, including waiting for a mutex $\text{mtx}(eid)$ and waiting for a duration tm etc.. ecbls is the name for the high-level abstract event control block (ECB) list β , which maps the event identifiers to abstract events, including abstract mutexes, abstract message queues and so on. The abstract mutex can be formalized as $\text{mutex}(pr, w)$, in which pr is the priority of the mutex and w is a pair of task identifiers and priorities. ecbls is used to implement mutexes. Note that the definition of the abstract state for ECB list is not tied to our framework, which means different kernel implementations may instantiate these abstract states with different settings.

Example of high-level specifications. We use $\mathbb{s}_{\text{dly}} \stackrel{\text{def}}{=} (\gamma_{\text{err}}(\text{ticks}) + (\gamma_{\text{dly}}(\text{ticks}); \text{sched}))$, as shown in Fig. 13, to specify the system API “void OSTimeDly(Int16u ticks)”, which delays the current task for the specified number of system ticks. The atomic operation $\gamma_{\text{err}}(\text{ticks})$ specifies the error case when $\text{ticks} = 0$. $\gamma_{\text{dly}}(\text{ticks})$ defines the atomic behavior of updating the status of the current task from “ready” to “waiting” with the duration set to ticks when $\text{ticks} > 0$, and the following **sched** switches to another ready task, following the scheduling policy specified by the abstract scheduler χ . Note that the exclusive conditions over ticks in $\gamma_{\text{err}}(\text{ticks})$ and $\gamma_{\text{dly}}(\text{ticks})$ make the non-deterministic choice statement deterministic. We omit the definitions of $\gamma_{\text{err}}(\text{ticks})$ and $\gamma_{\text{dly}}(\text{ticks})$ here.

As another example, below we show the abstract scheduler specification $\chi_{\mu\text{C}/\text{OS-II}}$ for $\mu\text{C}/\text{OS-II}$. It requires that the selected task be ready and have the highest priority among all the ready tasks.

$$\chi_{\mu\text{C}/\text{OS-II}} \stackrel{\text{def}}{=} \lambda \Sigma, t. \exists \alpha, pr. \Sigma(\text{tcbls}) = \alpha \wedge \alpha(t) = (pr, \text{rdy}) \wedge \forall t', pr'. (t \neq t' \wedge \alpha(t') = (pr', \text{rdy})) \rightarrow pr' \prec pr$$

The High-level Operational Semantics. Figure 12 shows selected rules of the high-level operational semantics, which are defined similarly as the low level. We use $\mathbb{P} \vdash \mathbb{W} \Rightarrow_H \mathbb{W}'$, $\mathbb{P} \vdash (C, \Delta, \Sigma) \xrightarrow{H} (C', \Delta, \Sigma')$ and $(s, \Sigma) \bullet \xrightarrow{H} (s', \Sigma')$ to represent the high-level program steps, the high-level task-steps, and the high-level kernel-steps, respectively. The operational semantics rules for program-steps correspond to an execution step of the abstract kernel specification, a step of the **sched** command, a step to handle events (abstractions for interrupts) and a step of client execution.

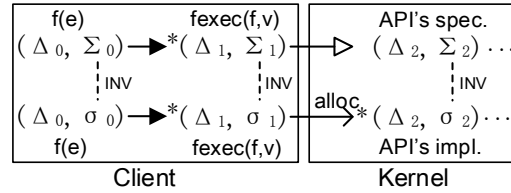


Fig. 14. Correspondence of Invoking APIs at two Levels

The high-level machine applies the API rule to execute the specification code of the function f . As shown in Fig.14, when the client invokes an API at two levels, both the low-level and high-level machines first execute with the exactly the same client-steps to the statement $\mathbf{fexec}(f, \bar{v})$. Then the low-level execution first allocates memory for arguments and local variables and enters the concrete function body of f as we explained before, while the high level enters the abstract specification code. The INAPI and ENDAPI rules are used to execute and return from the abstract function body. The EVT and ENDEVT rules are used for responding the external events (corresponding to low-level interrupts) at the high level. The HPRIM and ASSUME rules give the semantics for two basic abstract statements. The operational semantics of the sequential statements and non-deterministic choice statements are standard.

3.3 OS Correctness

As we explain in Sec.2.2, the correctness of OS kernels can be defined in terms of contextual refinement. Below we give its formal definition.

Definition 3.1 (OS Correctness). $O \sqsubseteq_{\psi} \mathbb{O}$ iff
 $\forall A, W, \mathbb{W}. \text{Match}(\psi, W, \mathbb{W}) \implies ((A, O), W) \preceq ((A, \mathbb{O}), \mathbb{W})$
 where $\psi \in \text{LOSFullSt} \rightarrow \text{HAbsSt} \rightarrow \text{Prop}$ and
 $\text{Match}(\psi, (T, \Delta, A, t), (T, \Delta, \Sigma)) \stackrel{\text{def}}{=} (t \in \text{dom}(T)) \wedge (\psi \wedge \Sigma) \wedge (t = \Sigma(\text{ctid})) \wedge (\text{dom}(T) = \text{dom}(\Sigma(\text{tcbls})))$

$$\begin{array}{c}
\text{(EvtTrace)} \quad \xi ::= \text{nil} \mid \downarrow \mid \varsigma :: \xi \\
\\
\frac{}{LETr(P, W, \epsilon)} \quad \frac{P \vdash W \Rightarrow_L \text{abort}}{LETr(P, W, \downarrow)} \quad \frac{P \vdash W \Rightarrow_L W' \quad LETr(P, W', \xi)}{LETr(P, W, \xi)} \\
\\
\frac{P \vdash W \Rightarrow_{\tilde{L}} W' \quad LETr(P, W', \xi)}{LETr(P, W, \varsigma :: \xi)} \\
\\
\frac{}{HETr(\mathbb{W}, \epsilon)} \quad \frac{\mathbb{P} \vdash \mathbb{W} \Rightarrow_H \text{abort}}{HETr(\mathbb{P}, \mathbb{W}, \downarrow)} \quad \frac{\mathbb{P} \vdash \mathbb{W} \Rightarrow_H \mathbb{W}' \quad HETr(\mathbb{P}, \mathbb{W}', \xi)}{HETr(\mathbb{P}, \mathbb{W}, \xi)} \\
\\
\frac{\mathbb{P} \vdash \mathbb{W} \Rightarrow_{\tilde{H}} \mathbb{W}' \quad HETr(\mathbb{P}, \mathbb{W}', \xi)}{HETr(\mathbb{P}, \mathbb{W}, \varsigma :: \xi)} \\
\\
(P, W) \preceq (\mathbb{P}, \mathbb{W}) \stackrel{\text{def}}{=} \forall \xi. LETr(P, W, \xi) \implies HETr(\mathbb{P}, \mathbb{W}, \xi)
\end{array}$$

Fig. 15. Event Trace Refinement

The low-level kernel code O refines its high-level abstract specifications \mathbb{O} with constraints ψ over initial kernel states, denoted as $O \sqsubseteq_{\psi} \mathbb{O}$, if and only if for any client code A , *low-level state* W and *high-level state* \mathbb{W} , if W and \mathbb{W} satisfy certain consistency constraint (w.r.t. ψ), then the set of observable behaviors of the low-level configuration $((A, O), W)$ is a subset of $((A, \mathbb{O}), \mathbb{W})$ (i.e., $(P, W) \preceq (\mathbb{P}, \mathbb{W})$, following the event trace refinement in [19]). The definition of $(P, W) \preceq (\mathbb{P}, \mathbb{W})$ is given in Fig. 15 and will be explained below.

The constraint **Match** requires that: (1) initially W and \mathbb{W} have the same task pool T and client state Δ ; (2) the current task t is in T ; (3) the low-level kernel state A and the high-level abstract state satisfy ψ ; (4) the *current* task at the low level and the high level are the same; and (5) the set of tasks in the abstract TCB list should be the same as those in the low-level task pool.

Event trace refinement for OS correctness. We give the definition of event trace refinement in Fig 15. nil means a empty trace. A trace is a sequence of externally observable events ς , and may end with a fault marker \downarrow . $LETr(P, W, \xi)$ means ξ can be generated by low-level machine (P, W) . $HETr(\mathbb{P}, \mathbb{W}, \xi)$ means ξ can be generated by high-level machine (\mathbb{P}, \mathbb{W}) . $(P, W) \preceq (\mathbb{P}, \mathbb{W})$ means that all the observable event trace generated by the low-level machine (P, W) can also be generated by high-level machine (\mathbb{P}, \mathbb{W}) .

4 Relational Program Logic for Refinement Verification

In this section, we present a CSL-style *relational* program logic for refinement verification. The logic uses relational assertions to prove refinement between an implementation and its specification. It also follows the ownership-transfer semantics in CSL to reason about multi-level hardware interrupts.

<pre> inc(){ int done=0, tmp; while(!done){ tmp=cnt; done=cas(&cnt,tmp,tmp+1) } } </pre>	<pre> {cnt = N} inc(); {cnt = N+1} </pre>	<pre> {∃N. cnt = N} inc(); {∃N. cnt = N} </pre>	<pre> {cnt = CNT ∧ [[⟨CNT++⟩]]} inc(); {cnt = CNT ∧ [[end]]} </pre>
(a) Implementation of <code>inc</code>	(b) Wrong spec.	(c) Weak spec.	(d) Refinement spec.

Fig. 16. Specification of Concurrent Programs

4.1 Refinement of concurrent programs, and relational reasoning.

For concurrent programs, refinement establishes stronger functional correctness than traditional Hoare triples. As an example, the function `inc` shown in Fig. 16(a) increments the counter `cnt`. It may be called simultaneously by concurrent tasks. Figure 16(b) gives pre-/post-conditions to specify `inc`, which would be valid in a sequential setting and is sufficient to describe the functionality. However, they cannot be used in a concurrent setting because they are not stable with respect to concurrent behaviors of other tasks. To make them stable, we may need the specifications in Fig. 16(c), which is too weak to capture the functionality.

Figure 16(d) gives a relational specifications to show that `inc` refines an abstract operation $\langle \text{CNT}++ \rangle$ [20], where $\langle C \rangle$ represents an *atomic* operation C . The relational assertions specify three important entities, the concrete state (`cnt`), the abstract state (`CNT`) and the abstract operation ($\langle \text{CNT}++ \rangle$) that the program refines (which could be non-atomic in general [20]). The precondition requires that initially `cnt` has the consistent value with its abstract counterpart `CNT`, and the abstract operation that `inc` needs to refine is $\langle \text{CNT}++ \rangle$. The postcondition ensures `cnt` and `CNT` remain consistent and the remaining abstract operation that needs to be refined is `end` (*i.e.*, $\langle \text{CNT}++ \rangle$ has been accomplished).

Our refinement proofs for OS kernels follow the same kind of relational reasoning, where the assertions now relate the concrete kernel state, the abstract kernel state (Σ) and the abstract statement (`s`).

4.2 Relational Assertion Language

Figure 17 gives the relational assertion language, and its semantics is given in Fig. 18.

$$\begin{aligned}
(Asrt) \quad p, q, r &::= \text{emp} \mid \text{empE} \mid x \mapsto v \mid \text{ISR}(isr) \mid \text{IE}(ie) \mid \text{IS}(is) \mid \text{CS}(cs) \mid \perp k \perp \mid \chi \triangleright t \\
&\quad \mid a \mapsto \Omega \mid [\![s]\!] \mid p * p \mid p \wedge p \mid \dots \\
(InvAsrt) \quad I &::= [p_0, \dots, p_N]
\end{aligned}$$

Fig. 17. Relational Assertions

As explained above, the assertions are interpreted over relational states Θ , which consist of the low-level task-local states σ , the high-level abstract states

Σ , and the abstract statements s that the low-level code needs to refine. Σ and s are defined in Fig. 10. σ , as shown in Fig. 18, consists of a task-local view m of program variables and memory, and also the global *isr* register and the task-local interrupt states δ (see Fig. 4). Here m contains the global and local variables (G and E respectively) and the memory M , whose definitions are omitted.

$$\begin{aligned}
(\text{RelState}) \quad \Theta &::= (\sigma, \Sigma, s) \\
(\sigma, \Sigma, s) \models \text{emp} &\quad \text{iff } \sigma.m.M = \emptyset \wedge \Sigma = \emptyset \\
(\sigma, \Sigma, s) \models \text{empE} &\quad \text{iff } \sigma.m.E = \emptyset \wedge (\sigma, \Sigma, s) \models \text{emp} \\
(\sigma, \Sigma, s) \models x \mapsto v &\quad \text{iff } \exists a. (\sigma.m.G)(x) = a \wedge \sigma.m.M = \{a \rightsquigarrow v\} \wedge \Sigma = \emptyset \\
(\sigma, \Sigma, s) \models \text{ISR}(isr') &\quad \text{iff } \sigma.isr = isr' \wedge (\sigma, \Sigma, s) \models \text{emp} \\
(\sigma, \Sigma, s) \models \perp k \perp &\quad \text{iff } ((k = N \wedge is = \text{nil}) \vee \exists is'. (\sigma.\delta.is = k :: is')) \wedge (\sigma, \Sigma, s) \models \text{emp} \\
(\sigma, \Sigma, s) \models \chi \triangleright t &\quad \text{iff } \chi \Sigma t \\
(\sigma, \Sigma, s) \models [\![s']\!] &\quad \text{iff } s = s' \wedge (\sigma, \Sigma, s) \models \text{emp} \\
(\sigma, \Sigma, s) \models a \mapsto \Omega &\quad \text{iff } \Sigma = \{a \rightsquigarrow \Omega\} \wedge \sigma.m.M = \emptyset \\
f \perp g &\stackrel{\text{def}}{=} \text{dom}(f) \cap \text{dom}(g) = \emptyset \quad \Sigma_1 \uplus \Sigma_2 \stackrel{\text{def}}{=} \begin{cases} \Sigma_1 \cup \Sigma_2 & \text{iff } \Sigma_1 \perp \Sigma_2 \\ \text{undef} & \text{otherwise} \end{cases} \\
\sigma_1 \uplus \sigma_2 &\stackrel{\text{def}}{=} \begin{cases} ((G, E, M_1 \cup M_2), isr, \delta) & \text{iff } M_1 \perp M_2 \wedge \sigma_1 = ((G, E, M_1), isr, \delta) \\ \text{undef} & \text{otherwise} \end{cases} \quad \wedge \sigma_2 = ((G, E, M_2), isr, \delta) \\
\Theta_1 \uplus \Theta_2 &\stackrel{\text{def}}{=} (\sigma_1 \uplus \sigma_2, \Sigma_1 \uplus \Sigma_2, s) \quad \text{where } \Theta_1 = (\sigma_1, \Sigma_1, s) \wedge \Theta_2 = (\sigma_2, \Sigma_2, s) \\
\Theta \models p_1 * p_2 &\quad \text{iff } \exists \Theta_1, \Theta_2. \Theta = \Theta_1 \uplus \Theta_2 \wedge \Theta_1 \models p_1 \wedge \Theta_2 \models p_2
\end{aligned}$$

Fig. 18. Semantics of Relational Assertions

Assertion **emp** says the low-level memory and the high-level abstract state are both empty. **empE** further requires that the local variable environment be empty too. $x \mapsto v$ specifies a singleton memory cell with v stored in the global program variable x . **ISR**(*isr*), **IS**(*is*), **IE**(*ie*) and **CS**(*cs*) specify the value of the corresponding interrupt status (see Fig. 4). $\perp k \perp$ means that the currently running interrupt handler is at level k (or $k = N$, meaning no running handlers).

$\chi \triangleright t$ says that, based on the high-level abstract state, the abstract scheduler χ picks t as the target task. $a \mapsto \Omega$ specifies a singleton high-level abstract state mapping the data name a to the abstract data Ω . $[\![s']\!]$ means the current abstract statement remaining to be refined is s . The separating conjunction $p_1 * p_2$ means p_1 and p_2 hold over disjoint parts of a relational state.

Ownership-transfer semantics for multi-level interrupts. CSL [22] prevents data races by enforcing disjoint ownership of resources among tasks. Synchronization is modeled in terms of ownership transfer. Feng *et al.* [11] extend CSL and assign ownership-transfer semantics to interrupt operations. The idea is demonstrated in Fig. 19, which shows the *logical* memory model when there are only one task and single-level interrupt. Since the interrupt handler can preempt the task, we let the handler to reserve its required memory first (represented as block B).

B must remain publicly available if the interrupt is enabled. Then the task can only access the remaining part (block T). We use grey boxes to represent *local* resources of the task. Disabling interrupts (**cli**) by the task essentially transfers the ownership of B from public to task-local. Correspondingly, **sti** converts the block from task-local to public, therefore the task *cannot* access it anymore. Similarly, invocation of the interrupt handler (not shown in the figure) automatically transfers B from public to the local resource of the handler, while **iret** transfers it back to public.

Since block B is shared between the interrupt handler and the task, it must be well-formed when it is public. We use the resource invariant I_0 to specify the well-formedness. Then the above ownership transfer semantics of **cli** and **sti** can be formalized in the following (simplified) program logic rules:

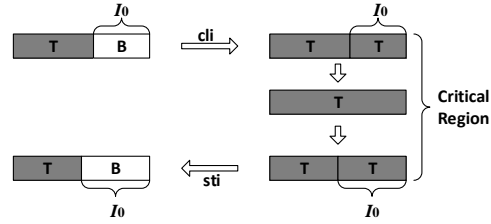


Fig. 19. Memory Partition for Handler and Non-Handler (Figure taken from [11])

$$\frac{}{I_0 \vdash \{p_t\} \mathbf{cli} \{p_t * I_0\}}$$

$$\frac{}{I_0 \vdash \{p_t * I_0\} \mathbf{sti} \{p_t\}}$$

Note that the partition between B and T is enforced *logically* using the separating conjunction in separation logic (see Fig. 18). It does not require physical separation in the program state model.

In this paper we extend this idea to support multi-level nested interrupts, where the ownership transfer of interrupt primitives is determined not only by the *ie* flag, but also by the *isr* register. Figure 20 shows the memory model (where the number N of interrupts is set to 6). Interrupt handlers at levels 0 to $N-1$ are assigned with resource blocks B_0, \dots, B_{N-1} respectively. B_N represents the resource shared only among tasks, *i.e.*, the non-handler code. We omit task-local resources, therefore there are *no* counterparts to block T in Fig. 19. Handlers' priorities to reserve their required resources are consistent with their interrupt priority levels. That is, B_0 satisfies all the need of the level-0 (highest priority) handler, while the level- k handler may need to access B_0, \dots, B_{k-1} , in addition to B_k . The non-handler has the lowest priority. Each block B_k is specified by the resource invariant $I(k)$, where I is defined as a sequence of $N+1$ assertions (see the assertion syntax defined above).

Figure 20 demonstrates the ownership transfer of resource caused by interrupt operations under different conditions. The grey or dotted blocks represent resources exclusively owned in interrupt handlers, different textures for different interrupts. The white ones represent resources *available* for share. Suppose initially we are at state (1), where the level-3 handler is being executed, as the value of *isr* indicates. Since interrupts are disabled, the handler owns $B_0 - B_3$, knowing *no* requests of levels 0 to 3 could be served. Enabling interrupts (**sti**) loses $B_0 - B_2$, as shown by state (2), but B_3 is remained because $isr(3) = 1$ and requests of the same (or lower) level are not handled. However, if $isr(3) = 0$

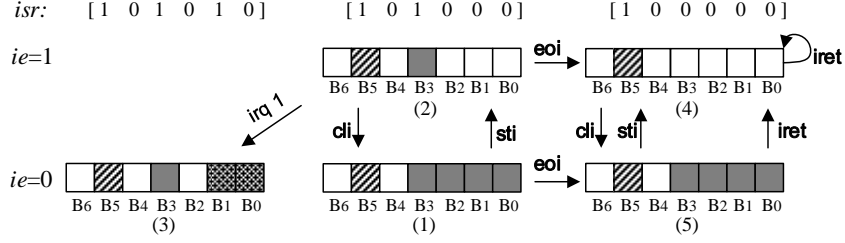


Fig. 20. Ownership-Transfer for Multi-level Interrupts

instead (as in state (5)), executing **sti** loses B_3 as well. Ownership transfer by **cli** is the dual of **sti**.

Executing **eoi** at state (1) leads to state (5), but it causes no ownership transfer because interrupts are disabled anyway. If interrupts are enabled instead, as in state (2), **eoi** loses the ownership of B_3 because another level-3 request may be handled in state (4). **ired** can be executed only after **eoi**. If interrupts are disabled (as in state (5)), it transfers $B_0 - B_3$ from local resources to shared resources. Otherwise (as in state (4)) there is no ownership transfer because the handler has lost the ownership of $B_0 - B_3$ already.

At state (2), interrupts with higher priority can be served. The “**irq 1**” step sets the bit $isr(1)$, disables interrupts, and transfers B_0 and B_1 from shared resources to local resources of the level-1 handler, as in state (3).

4.3 Inference Rules

The top rule. We show some selected program logic rules in Fig. 23. The TOPRULE establishes the judgment $\vdash_{\psi} O : \mathbb{Q}$, ensuring the correctness of O w.r.t. \mathbb{Q} if the initial concrete and abstract kernel states satisfy ψ (explained in Sec. 3.3).

$$\begin{aligned}
 (\text{FunPre}) \quad & fp \in \text{Vallist} \rightarrow \text{Asrt} & (\text{FunPost}) \quad & fq \in \text{Vallist} \rightarrow \text{Asrt} \\
 (\text{FunSpec}) \quad & \Gamma \in \text{FName} \rightarrow \text{FunPre} \times \text{FunPost}
 \end{aligned}$$

Fig. 21. Function Specifications

To verify the kernel, we need to come up with a specification Γ for the internal functions η_i in the low-level code, and a sequence of invariants I for kernel states. Γ defined in Fig. 21 assigns a pair of pre-/post-conditions to each internal function. The pre-/post-conditions is a mapping from value lists to assertions. The value list is used to specify the value of parameters.

Then we prove that the internal functions, the API implementations and the interrupt handlers in the low-level kernel satisfy their specifications, respectively (the last three premises in the first line of the TOPRULE rule). The proof of each component carries the abstract scheduler specification χ and the invariant I .

$$\begin{aligned}
(\sigma, \Sigma, \mathbf{s}) \models x \mapsto_l v & \quad \text{iff } \exists a. (\sigma.m.E)(x) = a \wedge \sigma.m.M = \{a \rightsquigarrow v\} \wedge \Sigma = \emptyset \\
\text{getD}(\eta, f) & \stackrel{\text{def}}{=} \begin{cases} \text{rev}(\mathcal{D}_1) ++ \mathcal{D}_2 & \text{iff } \eta(f) = (\tau, \mathcal{D}_1, \mathcal{D}_2, s) \\ \perp & \text{otherwise} \end{cases} \\
\text{BuildP}(\mathcal{D}, \bar{v}) & \stackrel{\text{def}}{=} \begin{cases} x \mapsto_l v * \text{BuildP}(\mathcal{D}', \bar{v}') & \text{iff } \mathcal{D} = (x, \tau) :: \mathcal{D}' \wedge \bar{v} = v :: \bar{v}' \\ x \mapsto_l - * \text{BuildP}(\mathcal{D}', \text{nil}) & \text{iff } \mathcal{D} = (x, \tau) :: \mathcal{D}' \wedge \bar{v} = \text{nil} \\ \text{emp} & \text{iff } \mathcal{D} = \text{nil} \\ \perp & \text{otherwise} \end{cases} \\
\text{BuildR}(\mathcal{D}) & \stackrel{\text{def}}{=} \text{BuildP}(\mathcal{D}, \text{nil}) \\
\text{BuildAPIPre}(\eta_a, f, \omega, \bar{v}) & \stackrel{\text{def}}{=} \text{OS}[\bar{0}, 1, \text{nil}, \text{nil}] * \text{BuildP}(\text{getD}(\eta_a, f), \bar{v}) * [\omega(\bar{v})] \\
\text{BuildAPIRet}(\eta_a, f) & \stackrel{\text{def}}{=} \text{OS}[\bar{0}, 1, \text{nil}, \text{nil}] * \text{BuildR}(\text{getD}(\eta_a, f)) * [\text{end}] \\
\text{BuildFunPre}(\eta_i, f, \bar{v}, fp) & \stackrel{\text{def}}{=} fp(\bar{v}) * \text{BuildP}(\text{getD}(\eta_i, f), \bar{v}) \\
\text{BuildFunRet}(\eta_i, f, \bar{v}, fq) & \stackrel{\text{def}}{=} fq(\bar{v}) * \text{BuildR}(\text{getD}(\eta_i, f)) \\
\text{BldltpPre}(k, \varepsilon, isr, is, I) & \stackrel{\text{def}}{=} \text{OS}[isr\{k \rightsquigarrow 1\}, 0, k :: is, \text{nil}] * I[0, k] * [\varepsilon(k)] * \text{empE} \\
\text{BldltpRet}(k, isr, is, I) & \stackrel{\text{def}}{=} \exists ie. \text{OS}[isr\{k \rightsquigarrow 0\}, ie, k :: is, \text{nil}] * \\
& \quad ((ie = 1 \wedge \text{emp}) \vee (ie = 0 \wedge I[0, k])) * [\text{end}] \\
I[n, m] & \stackrel{\text{def}}{=} \begin{cases} I(n) * I(n+1) * \dots * I(m) & \text{if } 0 \leq n \leq m \leq N \\ \text{emp} & \text{otherwise} \end{cases} \\
\text{OS}[isr, ie, is, cs] & \stackrel{\text{def}}{=} \exists k. \text{ISR}(isr) * \text{IE}(ie) * \text{IS}(is) * \text{CS}(cs) * \perp k \perp * \\
& \quad (\forall k'. 0 \leq k' < k \rightarrow isr(k') = 0) \\
[\psi] & \stackrel{\text{def}}{=} \lambda(\sigma, \Sigma, \mathbf{s}). \forall \Lambda. \psi \wedge \Sigma \wedge \exists t. \Lambda|_t = \sigma \\
\text{INV}(I, k) & \stackrel{\text{def}}{=} \exists isr. \text{ISR}(isr) * \\
& \quad ((isr(k) = 1 \wedge \text{emp}) \vee ((isr(k) = 0 \vee k = N) \wedge I(k))) \\
\text{SWINV}(I) & \stackrel{\text{def}}{=} \text{ISR}(\bar{0}) * \text{IE}(0) * (\exists k. \perp k \perp * I[0, k])
\end{aligned}$$

Fig. 22. Auxiliary Definitions of Inference Rules

The rule also requires that ψ ensures the initial states satisfy the invariant $I[0, N]$, the interrupt-related states are properly initialized, and the initial local variable environment is empty. $I[n, m]$ defined in Fig. 23 is the separating conjunction of invariants from level n to m . $\text{OS}[isr, ie, is, cs]$ specifies the status of interrupts, and requires that the currently executing handler (on top of is) have the highest priority among those in service (as recorded in isr). $[\psi]$ lifts ψ to relational assertions (defined in Fig. 22). More details about side conditions in the rule can be seen in Coq code [28].

Verifying interrupt handlers. The ITRP rule proves the correctness of interrupt handlers. It requires that each individual interrupt handler is correct with respect to its specification. The judgment for statements is in the form of $\Gamma; \chi; I; r; p_i \vdash \{p\} s \{q\}$. We follow the CSL-style reasoning, where I specifies shared resource blocks, and the pre-/post-conditions specify *local* resources that are accessed exclusively by the current task. The precondition is p , while q , r and p_i are all

$$\begin{array}{c}
\frac{O = (\eta_a, \eta_i, \theta) \quad \mathbb{O} = (\varphi, \varepsilon, \chi) \quad \chi; I \vdash \eta_i : \Gamma \quad \Gamma; \chi; I \vdash \eta_a : \varphi \quad \Gamma; \chi; I \vdash \theta : \varepsilon}{\vdash_\psi O : \mathbb{O}} \text{ (TOPRULE)} \\
\\
\frac{p = \text{BldltpPre}(k, \varepsilon, isr, is, I) \quad p_i = \text{BldltpRet}(k, isr, is, I) \quad \text{dom}(\theta) = \text{dom}(\varepsilon) \quad \Gamma; \chi; I; \text{false}; p_i \vdash \{p\} \theta(k) \{ \text{false} \} \quad \text{for all } k \in \{0, \dots, N-1\}}{\Gamma; \chi; I \vdash \theta : \varepsilon} \text{ (ITRP)} \\
\\
\frac{\text{dom}(\eta) = \text{dom}(\varphi) \quad \text{BuildAPIPre}(\eta, f, \omega, \bar{v}) = p \quad \text{BuildAPIRet}(\eta, f) = r \quad \eta(f) = (-, -, -, s) \quad \varphi(f) = \omega \quad \Gamma; \chi; I; r; \text{false} \vdash \{p\} s \{ \text{false} \}}{\Gamma; \chi; I \vdash \eta : \varphi} \text{ (WF API)} \\
\\
\frac{\text{dom}(\eta) = \text{dom}(\Gamma) \quad \text{BuildFunPre}(\eta, f, \bar{v}, fp) = p \quad \text{BuildFunRet}(\eta, f, \bar{v}, fq) = r \quad \eta(f) = (-, -, -, s) \quad \Gamma(f) = (fp, fq) \quad \Gamma; \chi; I; r; \text{false} \vdash \{p\} s \{ \text{false} \}}{\chi; I \vdash \eta : \Gamma} \text{ (WFFUN)} \\
\\
\hline
\\
\frac{p \Rightarrow p' \quad \Gamma; \chi; I; r; p_i \vdash \{p'\} s \{q'\} \quad q' \Rightarrow q}{\Gamma; \chi; I; r; p_i \vdash \{p\} s \{q\}} \text{ (CONSEQ)} \\
\\
\frac{\Gamma; \chi; I; r; p_i \vdash \{p_1\} s_1 \{p_2\} \quad \Gamma; \chi; I; r; p_i \vdash \{p_2\} s_2 \{p_3\}}{\Gamma; \chi; I; r; p_i \vdash \{p_1\} s_1; s_2 \{p_3\}} \text{ (SEQ)} \\
\\
\frac{\Gamma; \chi; I; r; p_i \vdash \{p_1\} s \{p_2\} \quad q \text{ does not specify } ie, is, cs, isr \text{ and } s}{\Gamma; \chi; I; r * q; p_i * q \vdash \{p_1 * q\} s \{p_2 * q\}} \text{ (FRM)} \\
\\
\frac{p \Rightarrow p' \quad \Gamma; \chi; I; r; p_i \vdash \{p'\} s \{q'\} \quad q' \Rightarrow q}{\Gamma; \chi; I; r; p_i \vdash \{p\} s \{q\}} \text{ (ABSCSQ)} \\
\\
\hline
\\
\frac{}{\Gamma; \chi; I; r; p_i \vdash \{ \text{OS}[isr, 1, is, cs] * \perp k \perp * [\mathbb{S}] \} \text{enct} \{ \text{OS}[isr, 0, is, 1 :: cs] * \text{INV}(I, k) * I[0, k-1] * [\mathbb{S}] \}} \text{ (ENCRT)} \\
\\
\frac{}{\Gamma; \chi; I; r; p_i \vdash \{ \text{OS}[isr, 0, is, cs] * [\mathbb{S}] \} \text{enct} \{ \text{OS}[isr, 0, is, 0 :: cs] * [\mathbb{S}] \}} \text{ (ENCRT-0)} \\
\\
\frac{}{\Gamma; \chi; I; r; p_i \vdash \{ \text{OS}[isr, 0, is, 1 :: cs] * \perp k \perp * \text{INV}(I, k) * I[0, k-1] * [\mathbb{S}] \} \text{exct} \{ \text{OS}[isr, 1, is, cs] * [\mathbb{S}] \}} \text{ (EXCRT)} \\
\\
\frac{}{\Gamma; \chi; I; r; p_i \vdash \{ \text{OS}[isr, 0, is, 0 :: cs] * [\mathbb{S}] \} \text{exct} \{ \text{OS}[isr, 0, is, cs] * [\mathbb{S}] \}} \text{ (EXCRT-0)} \\
\\
\frac{}{\Gamma; \chi; I; r; p_i \vdash \{ \text{OS}[isr, 1, k :: is, cs] * I(k) * [\mathbb{S}] \} \text{eoi } k \{ \text{OS}[isr\{k \rightsquigarrow 0\}, 1, k :: is, cs] * [\mathbb{S}] \}} \text{ (EOI)} \\
\\
\frac{p \Leftrightarrow \text{SWINV}(I) * \text{IS}(is) * \text{CS}(cs)}{\Gamma; \chi; I; r; p_i \vdash \{ (p * [\text{sched}; \mathbb{S}]) \wedge \chi \triangleright x \} \text{switch } x \{ p * [\mathbb{S}] \}} \text{ (SWITCH)} \\
\\
\frac{p \Rightarrow p_i}{\Gamma; \chi; I; \text{false}; p_i \vdash \{p\} \text{iext} \{ \text{false} \}} \text{ (IEXT)} \quad \frac{p \Rightarrow r}{\Gamma; \chi; I; r; \text{false} \vdash \{p\} \text{return} \{ \text{false} \}} \text{ (RET)}
\end{array}$$

Fig. 23. Selected Inference Rules

post-conditions for different exits, *i.e.*, sequential composition, return from functions, and return from interrupts, respectively. For the whole body of interrupt handlers, we disable the other two exits by setting r and q to **false**.

We build the pre-/post-conditions of interrupt handlers with the auxiliary definitions **BldltpPre** and **BldltpRet** given in Fig. 22. The precondition says that, when entering the level- k handler, $isr(k)$ is set to 1, the interrupt is disabled and k is pushed onto the interrupt stack is (therefore $OS[isr\{k \rightsquigarrow 1\}, 0, k :: is, nil]$). Since there is no handler of higher-priority in service, the handler has exclusive access to the resource $I[0, k]$ (see Fig. 20). It also needs to refine the high-level specification code $\varepsilon(k)$. **empE** requires there are no local variables at the beginning. The built post-condition requires that: (1) the corresponding isr bit has been cleared; (2) if interrupts are enabled ($ie = 1$), the handler has no access to the shared resources; otherwise it needs to ensure that its owned resources are well formed w.r.t. $I[0, k]$ (see the two **iret** steps in Fig. 20); and (3) there is no high-level specification code remaining to be refined (*i.e.*, the abstract specification code $\varepsilon(k)$ specified in the precondition has been fulfilled).

Similarly, we use **BuildAPIPre** and **BuildAPIRet** defined in Fig. 22 to construct the pre-/post-conditions for kernel APIs. The local states before and after calling to the API f are specified by **BuildP**($getD(\eta_a, f), \bar{v}$) and **BuildAPIRet**(η_a, f), which specify the memory locations of the arguments and local variables of the API. For the internal functions, in addition to these local states, we need to add the local states specified by the functions specifications. The rules of proving $\chi; I \vdash \eta_i : F$ and $F; \chi; I \vdash \eta_a : \varphi$ are similar to the rules for interrupt handlers.

In the middle of Fig.23, we give the **SEQ** rule for sequential compositionality, and **CONSEQ** rule to strengthen and weaken the precondition and postcondition respectively. Also as in separation logic, the **FRM** is designed for modular reasoning, in which the side-condition requires that the framed assertion should not say anything about interrupt states and abstract statements.



Rules for commands. The **IEXT** rule simply requires that the post-condition p_i holds when we reach the end of the interrupt handler. The **RET** rule requires that the post-condition r holds when we reach the end of the non-handler function. The **ENCRT** rule shows the ownership transfer when interrupts are disabled. Suppose we are at the level- k handler ($k = N$ means we are executing the non-handler code). Disabling interrupts prevents interrupt requests from level 0 to $k - 1$, therefore the current task gains the ownership of $I[0, k - 1]$. The transfer of the k -th block is specified by **INV**(I, k) in Fig. 23. If the bit $isr(k)$ is 0 (or $k = N$), the task also gains the ownership of $I(k)$, otherwise it already has the ownership of the k -th block and there is no extra ownership transfer. The two scenarios are also demonstrated by the two **cli** steps in Fig. 20. If interrupts are already disabled when **encrt** is executed, there is no ownership transfer, as shown by the **ENCRT-0** rule.

The **EXCRT** rule is the dual of the **ENCRT** rule (see the two **sti** steps in Fig. 20). Correspondingly there is a **EXCRT-0** rule. The **EOI** rule says, if interrupts are enabled, the task loses the ownership of $I(k)$ after **eo** k . Otherwise there is no

ownership transfer and the corresponding rule is omitted (see the two **eo**i steps in Fig. 20).

The **SWITCH** rule requires that the invariant $\text{SWINV}(I)$ holds before switching away and it is preserved after switching back. $\text{SWINV}(I)$, defined in Fig. 23, says that interrupts must be disabled, and all the bits of *isr* are 0 (*i.e.*, either we are running non-handler code or we are in the outmost layer of nested invocation of interrupt handlers and have already executed **eo**i). Also if we are running level- k code (either handler or non-handler if $k = N$), the resource blocks 0 to k acquired before should satisfy $I[0, k]$, so that the target task could access them. The rule also says that the task-local states *is* and *cs* are not changed by **switch**.

To establish refinement, the precondition also requires that the high-level abstract scheduler χ picks the same task with the one in x , and **switch** x at the low level correspond to the **sched** step at the high level. Therefore in the post-condition **sched** is no longer in the remaining abstract operations.

Following [20], the ABSCSQ rule looks like a regular consequence rule but allows us to *execute* the abstract code. The implication $p \Rightarrow p'$ is defined below.

$$\forall \sigma, \Sigma, s. ((\sigma, \Sigma, s) \models p) \implies \exists \Sigma', s'. ((s, \Sigma) \bullet \text{H} \dashrightarrow^* (s', \Sigma')) \wedge ((\sigma, \Sigma', s') \models p')$$

That is, given a related state (σ, Σ, s) satisfying p , the abstract code s could execute zero or multiple steps starting from Σ and reach (Σ', s') , so that the resulting related state (σ, Σ', s') satisfies p' . This rule allows us to establish simulation between the concrete and the abstract code, which then ensures refinement.

We can look at Fig. 16 to see the use of this rule. Suppose we want to verify **inc**() using the specification in Fig. 16(d). When we reach the **cas** command (see Fig. 16(a)), we have the precondition $(\text{tmp} = \text{cnt} \wedge \text{cnt} = \text{CNT} \wedge [\langle \text{CNT}++ \rangle] \vee \dots)$ (the case for $\text{tmp} \neq \text{cnt}$ omitted). Right after **cas**, we have $(\text{done} \wedge \text{cnt} = \text{CNT} + 1 \wedge [\langle \text{CNT}++ \rangle] \vee \neg \text{done} \wedge \dots)$. We have $\text{cnt} = \text{CNT} + 1$ because **cnt** increments if **cas** succeeds. To establish the simulation, we apply the ABSCSQ rule to execute the abstract code, because $(\text{cnt} = \text{CNT} + 1 \wedge [\langle \text{CNT}++ \rangle]) \Rightarrow (\text{cnt} = \text{CNT} \wedge [\langle \text{end} \rangle])$, following the above definition of $p \Rightarrow p'$.

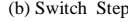
Theorem 4.12 gives the soundness of the framework. The proofs are based on a compositional simulation following [19], and have been formalized in Coq. More details about the logic can be seen in TR [28].

Theorem 4.1 (Soundness). $\vdash_\psi O : \mathbb{O} \implies O \sqsubseteq_\psi \mathbb{O}$.

4.4 Soundness via. Simulations

In this section, we give the semantics of our logic judgments and show the proof sketch of proving the soundness of the CSL-style relational logic.

Judgment Semantics. In this section, we show the semantics of inference rules. The semantics of $I; \chi; I; r; p_i \vdash \{p\} s \{q\}$ is defined in Def. 4.18 via a compositional simulation (defined in Def 4.17). Def. 4.4, 4.5 and 4.6 are the semantics of **WFINT** rule, **WFAPI** rule and **WFFUN** rule, respectively.



For each step in the kernel method, we require it to be safe.

$$\begin{aligned}
\sigma \perp \sigma' &\stackrel{\text{def}}{=} \sigma = ((G, E, M), \text{isr}, \delta) \wedge \sigma = ((G, E, M'), \text{isr}, \delta) \wedge M \perp M' \\
I\{n, m\} &\stackrel{\text{def}}{=} \begin{cases} \text{INV}(I, n) * \text{INV}(I, n+1) * \dots * \text{INV}(I, m) & \text{if } 0 \leq n \leq m \leq N \\ \text{emp} & \text{otherwise} \end{cases} \\
\lfloor I \rfloor &\stackrel{\text{def}}{=} ((\text{IE}(1) * I\{0, N\}) \vee (\text{IE}(0) * (\exists k. \perp k \perp * I\{k+1, N\})))
\end{aligned}$$

Fig. 25. Auxiliary Definitions of Simulations

Definition 4.2 (Judgment Semantics). $\Gamma; \chi; I; r; p_i \models \{p\} s \{q\}$ holds, iff for any σ, Σ and \mathfrak{s} , if $(\sigma, \Sigma, \mathfrak{s}) \models p$, then $\Gamma; \chi; I; r; p_i; q \models ((\sigma, (\circ, \bullet)), \sigma) \preceq (\mathfrak{s}, \Sigma)$.

Definition 4.3 (Method Simulation). $\Gamma; \chi; I; r; p_i; q \models (C, \sigma) \preceq (\mathfrak{s}, \Sigma)$ holds, whenever:

- **Normal Steps:** for any $P, C', \sigma_s, \Sigma_s, \sigma_1$ and σ'_1 , if $C \neq (\mathbf{fexec}(_, _), _)$, $(\sigma_s, \Sigma_s, _) \models \lfloor I \rfloor$, $\sigma_1 = \sigma \uplus \sigma_s$, $\Sigma \perp \Sigma_s$ and $P \vdash (C, \sigma_1) \bullet \text{L} \rightarrow (C', \sigma'_1)$, then there exist $\Sigma'_s, \mathfrak{s}', \Sigma', \sigma'$ and σ'_s , such that the followings hold:
 - $\sigma'_1 = \sigma' \uplus \sigma'_s$, $(\sigma'_s, \Sigma'_s, _) \models \lfloor I \rfloor$,
 - $(\mathfrak{s}, \Sigma \uplus \Sigma_s) \bullet \text{H} \rightarrow^* (\mathfrak{s}', \Sigma' \uplus \Sigma'_s)$,
 - $\Gamma; \chi; I; r; p_i; q \models (C', \sigma') \preceq (\mathfrak{s}', \Sigma')$.
- **Function Call:** for any $\sigma_s, \Sigma_s, \kappa_s, f$ and \bar{v} , if $C = (\mathbf{fexec}(f, \bar{v}), (\circ, \kappa_s))$, $\sigma \perp \sigma_s$, $(\sigma_s, \Sigma_s, _) \models \lfloor I \rfloor$ and $\Sigma \perp \Sigma_s$, then there exist $\sigma_1, \sigma_f, \Sigma_1, \Sigma_f, \Sigma'_s, \Sigma'_f, \mathfrak{s}', fp$ and fq , such that the followings hold:
 - $\Gamma(f) = (fp, fq)$,
 - $(\mathfrak{s}, \Sigma \uplus \Sigma_s) \bullet \text{H} \rightarrow^* (\mathfrak{s}', \Sigma' \uplus \Sigma'_s)$, $(\sigma_s, \Sigma'_s, \mathfrak{s}') \models \lfloor I \rfloor$,
 - $\sigma = \sigma_1 \uplus \sigma_f$, $\Sigma' = \Sigma_1 \uplus \Sigma_f$, $(\sigma_1, \Sigma_1, \mathfrak{s}') \models fp(\text{rev}(\bar{v}))$,
 - for any $\sigma', \sigma'_1, \Sigma'', \Sigma'_1$ and \mathfrak{s}'' , if $\sigma.m.G = \sigma'.m.G$, $\sigma.m.E = \sigma'.m.E$, $(\sigma'_1, \Sigma'_1, \mathfrak{s}'') \models fq(\text{rev}(\bar{v}))$, $\sigma' = \sigma'_1 \uplus \sigma_f$, and $\Sigma'' = \Sigma'_1 \uplus \Sigma_f$, then $\Gamma; \chi; I; r; p_i; q \models ((\mathbf{skip}, (\circ, \kappa_s)), \sigma') \preceq (\mathfrak{s}'', \Sigma'')$.
- **Context Switch:** for any σ_s, κ_s, x and Σ_s , if $\sigma \perp \sigma_s$, $C = (\mathbf{switch} \ x, (\circ, \kappa_s))$, $(\sigma_s, \Sigma_s, _) \models \lfloor I \rfloor$ and $\Sigma \perp \Sigma_s$, then there exist $\sigma_1, \sigma', \Sigma_1, \Sigma', \Sigma'_s$ and \mathfrak{s}' , such that the followings hold:
 - $(\mathfrak{s}, \Sigma \uplus \Sigma_s) \bullet \text{H} \rightarrow^* (\mathbf{sched}; \mathfrak{s}', \Sigma_1 \uplus \Sigma' \uplus \Sigma'_s)$,
 - $(\sigma_s, \Sigma'_s, _) \models \lfloor I \rfloor$, $\sigma = \sigma_1 \uplus \sigma'$,
 - $(\sigma', \Sigma', _) \models \text{SWINV}(I) \wedge (\chi \triangleright x)$
 - for any $\sigma'', \sigma''', \Sigma''$ and Σ''' , if $\sigma''' = \sigma_1 \uplus \sigma''$, $\Sigma''' = \Sigma_1 \uplus \Sigma''$ and $(\sigma'', \Sigma'', _) \models \text{SWINV}(I)$, then $\Gamma; \chi; I; r; p_i; q \models ((\mathbf{skip}, (\circ, \kappa_s)), \sigma''') \preceq (\mathfrak{s}', \Sigma''')$.
- **Skip:** for any σ_s and Σ_s , if $C = (\mathbf{skip}, (\circ, \bullet))$, $\sigma \perp \sigma_s$, $(\sigma_s, \Sigma_s, _) \models \lfloor I \rfloor$ and $\Sigma \perp \Sigma_s$, then there exist Σ', Σ'_s and \mathfrak{s}' , such that $(\mathfrak{s}, \Sigma \uplus \Sigma_s) \bullet \text{H} \rightarrow^* (\mathfrak{s}', \Sigma' \uplus \Sigma'_s)$, $(\sigma_s, \Sigma'_s, _) \models \lfloor I \rfloor$ and $(\sigma, \Sigma', \mathfrak{s}') \models q$;
- **IRet:** for any κ_s, σ_s and Σ_s , if $C = (\mathbf{ixext}, (\circ, \kappa_s))$, $\lfloor \kappa_s \rfloor = \perp$, $\lfloor \kappa_s \rfloor_c = \perp$, $\sigma \perp \sigma_s$, $(\sigma_s, \Sigma_s, _) \models \lfloor I \rfloor$ and $\Sigma \perp \Sigma_s$, then there exist Σ', Σ'_s and \mathfrak{s}' , such that $(\mathfrak{s}, \Sigma \uplus \Sigma_s) \bullet \text{H} \rightarrow^* (\mathfrak{s}', \Sigma' \uplus \Sigma'_s)$, $(\sigma_s, \Sigma'_s, _) \models \lfloor I \rfloor$, and $(\sigma, \Sigma', \mathfrak{s}') \models p_i$;

- **ReturnE** and **Return** cases are similar to the **IRet** case, we omit them here;
- **Abort**: for any P, Σ_s, σ_s and σ' , if $C \neq (\mathbf{fexec}(-, -), -)$, $\sigma' = \sigma \uplus \sigma_s$, $(\sigma_s, \Sigma_s, -) \models [I]$ and $\Sigma \perp \Sigma_s$, then $\neg(P \vdash (C, \sigma') \bullet \text{---} \text{abort})$.

Definition 4.4 (Well-Formed Interrupts). $\Gamma; \chi; I \models \theta : \varepsilon$ holds, iff for any k, isr, is, G, p and p_i , if $\varepsilon(k) = \mathbb{s}$, $p = \text{BldltpPre}(k, \mathbb{s}, isr, is, I)$, $p_i = \text{BldltpRet}(k, isr, is, I)$, then there exists s , such that $\theta(k) = s$ and $\Gamma; \chi; I; \text{false}; p_i \models \{p\}s\{\text{false}\}$

Definition 4.5 (Well-Formed APIs). $\Gamma; \chi; I \models \eta_a : \varphi$ holds, iff for any f, \bar{v}, ω, p and r , if $\varphi(f) = \omega$, $p = \text{BuildAPIPre}(\eta_a, f, \omega, \bar{v})$, $r = \text{BuildAPIRet}(\eta_a, f)$, then there exists s , such that $\eta_a(f) = (\tau, \mathcal{D}_1, \mathcal{D}_2, s)$ and $\Gamma; \chi; I; r; \text{false} \models \{p\}s\{\text{false}\}$.

Definition 4.6 (Well-Formed Internal Functions). $\chi; I \models \eta_i : \Gamma$ holds, iff $\text{dom}(\Gamma) = \text{dom}(\eta_i)$ and for any f, fp, fq, \bar{v}, p and r , if $\Gamma(f) = (fp, fq)$, $p = \text{BuildFunPre}(\eta_i, f, \bar{v}, fp)$ and $r = \text{BuildFunRet}(\eta_i, f, \bar{v}, fq)$ then there exists s , such that $\eta_i(f) = (-, -, -, s)$ and $\Gamma; \chi; I; r; \text{false} \models \{p\}s\{\text{false}\}$.

Soundness Proof. To prove the soundness of our logic (Theorem 4.12), following [18], we need to define the following two simulation relations as bridges:

1. A whole-program simulation relation (defined in Def. 4.7) between concrete and abstract levels, which implies that the observable behaviors of the low level is a subset of the high level (see Lemma 4.9);
2. A task-local simulation relations (defined in Def. 4.8) between the low-level task and the high-level task, which can be composed together to ensure the whole program simulation (see Lemma 4.10). And it could be implied by the method simulation with some side conditions (see Lemma 4.11).

Definition 4.7 (Program Simulation). $(P, W) \preceq (\mathbb{P}, \mathbb{W})$ holds, whenever:

- for any W, W' and P , if $P \vdash W =_{L \Rightarrow} W'$, then there exist \mathbb{W}' , such that the followings hold:
 - $\mathbb{P} \vdash \mathbb{W} =_{H \Rightarrow^*} \mathbb{W}'$ and $W'.\Delta = \mathbb{W}'.\Delta$,
 - $(P, W') \preceq (\mathbb{P}, \mathbb{W}')$.
- for any W, W' and P , if $P \vdash W =_{\tilde{L} \Rightarrow} W'$, then there exist \mathbb{W}' , such that the followings hold:
 - $\mathbb{P} \vdash \mathbb{W} =_{\tilde{H} \Rightarrow^*} \mathbb{W}'$ and $W'.\Delta = \mathbb{W}'.\Delta$,
 - $(P, W') \preceq (\mathbb{P}, \mathbb{W}')$.
- for any W, W' and P , if $P \vdash W =_{L \Rightarrow} \text{abort}$, then $\mathbb{P} \vdash \mathbb{W} =_{H \Rightarrow^*} \text{abort}$

Definition 4.8 (Task Simulation). $P; \mathbb{P}; I; p \models (C_l, \sigma) \preceq (C_h, \Sigma)$ holds, whenever:

- **Normal Steps:** for any $C'_l, \Delta, \Delta', \sigma_s, \Sigma_s, \sigma_1$ and σ'_1 , if $(\sigma_s, \Sigma_s, -) \models [I]$, $\sigma_1 = \sigma \uplus \sigma_s$, $\Sigma \perp \Sigma_s$ and $P \vdash (C, \Delta, \sigma_1) \xrightarrow{L} (C', \Delta', \sigma'_1)$, then there exist $\Sigma'_s, C'_h, \Sigma', \sigma'$ and σ'_s , such that the followings hold:
 - $\sigma'_1 = \sigma' \uplus \sigma'_s$, $(\sigma'_s, \Sigma'_s, -) \models [I]$,
 - $\mathbb{P} \vdash (C_h, \Delta, \Sigma \uplus \Sigma_s) \xrightarrow{H}^* (C'_h, \Delta', \Sigma' \uplus \Sigma'_s)$,
 - $P; \mathbb{P}; I; p \models (C'_l, \sigma') \preceq (C'_h, \Sigma')$.
- **Event Steps:** for any $C'_l, \Delta, \Delta', \sigma_s, \Sigma_s, \sigma_1$ and σ'_1 , if $(\sigma_s, \Sigma_s, -) \models [I]$, $\sigma_1 = \sigma \uplus \sigma_s$, $\Sigma \perp \Sigma_s$ and $P \vdash (C, \Delta, \sigma_1) \xrightarrow{\tilde{L}} (C', \Delta', \sigma'_1)$, then there exist $\Sigma'_s, C'_h, \Sigma', \sigma'$ and σ'_s , such that the followings hold:
 - $\sigma'_1 = \sigma' \uplus \sigma'_s$, $(\sigma'_s, \Sigma'_s, -) \models [I]$,
 - $\mathbb{P} \vdash (C_h, \Delta, \Sigma \uplus \Sigma_s) \xrightarrow{\tilde{H}}^* (C'_h, \Delta', \Sigma' \uplus \Sigma'_s)$,
 - $P; \mathbb{P}; I; p \models (C'_l, \sigma') \preceq (C'_h, \Sigma')$.
- **Context Switch:** for any $\Delta, \sigma_s, \kappa_s, x$ and Σ_s , if $\sigma \perp \sigma_s$, $C = (\mathbf{switch} \ x, (\circ, \kappa_s))$, $(\sigma_s, \Sigma_s, -) \models [I]$ and $\Sigma \perp \Sigma_s$, then there exist $\sigma_1, \sigma', \Sigma', \Sigma_1, \Sigma'_s, \chi, \mathbb{s}'$ and K , such that the followings hold:
 - $\mathbb{P} = (-, (-, -, \chi))$
 - $\mathbb{P} \vdash (C_h, \Delta, \Sigma \uplus \Sigma_s) \xrightarrow{H}^* ((\mathbf{sched}; \mathbb{s}', K), \Delta, \Sigma_1 \uplus \Sigma' \uplus \Sigma'_s)$,
 - $(\sigma_s, \Sigma'_s, -) \models [I]$, $\sigma = \sigma_1 \uplus \sigma'$,
 - $(\sigma', \Sigma', -) \models \mathbf{SWINV}(I) \wedge (\chi \triangleright x)$
 - for any $\sigma'', \sigma''', \Sigma''$ and Σ''' , if $\sigma''' = \sigma_1 \uplus \sigma''$, $\Sigma''' = \Sigma_1 \uplus \Sigma''$ and $(\sigma'', \Sigma'', -) \models \mathbf{SWINV}(I)$, then $P; \mathbb{P}; I; p \models ((\mathbf{skip}, (\circ, \kappa_s)), \sigma''') \preceq ((\mathbb{s}', K), \Sigma''')$.
- **Skip:** for any Δ, σ_s and Σ_s , if $C = (\mathbf{skip}, (\circ, \bullet))$, $\sigma \perp \sigma_s$, $(\sigma_s, \Sigma_s, -) \models [I]$ and $\Sigma \perp \Sigma_s$, then there exist Σ', Σ'_s , such that $(\sigma_s, \Sigma'_s, -) \models [I]$, $(\sigma, \Sigma', -) \models p$ and $\mathbb{P} \vdash (C_h, \Delta, \Sigma \uplus \Sigma_s) \xrightarrow{H}^* ((\mathbf{skip}, (\circ, \bullet)), \Delta, \Sigma' \uplus \Sigma'_s)$.
- **Abort:** for any $\Delta, \Sigma_s, \sigma_s$ and σ' , if $\sigma' = \sigma \uplus \sigma_s$, $(\sigma_s, \Sigma_s, -) \models [I]$ and $\Sigma \perp \Sigma_s$ and $P \vdash (C_l, \Delta, \sigma') \xrightarrow{L} \mathbf{abort}$, then $\mathbb{P} \vdash (C_h, \Delta, \Sigma \uplus \Sigma_s) \xrightarrow{H}^* \mathbf{abort}$.

Lemma 4.9 (ProgSim Implies Event Trace Refinement). For any P, \mathbb{P}, W and \mathbb{W} , if $(P, W) \preceq (\mathbb{P}, \mathbb{W})$, then $(P, W) \preceq (\mathbb{P}, \mathbb{W})$.

Lemma 4.10 (Compositionality). For any $\eta_a, \eta_i, \theta, \varphi, \varepsilon, \chi, \psi, T_l, T_h, \Delta, t_c$ and Σ , if the followings holds:

- $P = (A, (\eta_a, \eta_i, \theta))$, $\mathbb{P} = (A, (\varphi, \varepsilon, \chi))$
- $\mathbf{Match}(\psi, (T_l, \Lambda, \Delta, t_c), (T_h, \Sigma, \Delta))$, $\Lambda = ((G, \Pi, M), \text{isr}, \pi)$
- $\Gamma; \chi; I \models \theta : \varepsilon$, $\chi; I \models \eta_i : \Gamma$, $\chi; I \models \eta_i : \Gamma$
- $T_l = \{t_1 \rightsquigarrow C_{l1}, \dots, t_n \rightsquigarrow C_{ln}\}$, $T_h = \{t_1 \rightsquigarrow C_{h1}, \dots, t_n \rightsquigarrow C_{hn}\}$
- $M = M_1 \uplus M_2 \uplus \dots \uplus M_n \uplus M_s$, $\Sigma = \Sigma_1 \uplus \Sigma_2 \uplus \dots \uplus \Sigma_n \uplus \Sigma_s$
- $\Lambda|_{t_c} = \sigma_c$, $(\sigma_c \triangleleft M_s, \Sigma_s, -) \models [I]$
- $P; \mathbb{P}; I; p \models (C_{lc}, \sigma_c \triangleleft M_c) \preceq (C_{hc}, \Sigma_c)$
- for any $i, \sigma_i, \sigma_r, \Sigma_r$, $i \neq c$, $\sigma_i = (\Lambda|_{t_i}) \triangleleft M_i$, $(\sigma_r, \Sigma_r, -) \models \mathbf{SWINV}(I)$, $\sigma_i \perp \sigma_r$, $\Sigma_i \perp \Sigma_r$, then $P; \mathbb{P}; I; p \models (C_{li}, \sigma_i \uplus \sigma_r) \preceq (C_{hi}, \Sigma_i \uplus \Sigma_r)$

then $(P, (T_l, \Lambda, \Delta, t_c)) \preceq (\mathbb{P}, (T_h, \Sigma, \Delta))$.

Lemma 4.11. For any s, P, \mathbb{P}, I ,

- $P = (A, (\eta_a, \eta_i, \theta)), \mathbb{P} = (A, (\varphi, \varepsilon, \chi))$
- $(\sigma, \Sigma, _) \models \text{OS}[\bar{0}, 1, \text{nil}, \text{nil}]$
- $\Gamma; \chi; I \models \theta : \varepsilon, \quad \chi; I \models \eta_i : \Gamma, \quad \chi; I \models \eta_i : \Gamma$

then $P; \mathbb{P}; I; \text{OS}[\bar{0}, 1, \text{nil}, \text{nil}] \models ((s, (\circ, \bullet)), \sigma) \preceq ((s, (\circ, \bullet)), \Sigma)$.

Here we present the proof sketch of Theorem 4.12, and omit the proofs of other lemmas. The complete proofs can be found in our Coq proof [?]. Note the co-inductive proofs done in Coq for the soundness theorem are non-trivial.

Theorem 4.12 (Soundness).

1. $\Gamma; \chi; I; r; p_i \vdash \{p\} s \{q\} \implies \Gamma; \chi; I; r; p_i \models \{p\} s \{q\}$
2. $\Gamma; \chi; I \vdash \theta : \varepsilon \implies \Gamma; \chi; I \models \theta : \varepsilon$
3. $\Gamma; \chi; I \vdash \eta_a : \varphi \implies \Gamma; \chi; I \models \eta_a : \varphi$
4. $\chi; I \vdash \eta : \Gamma \implies \chi; I \models \eta : \Gamma$
5. $\vdash_\psi O : \mathbb{O} \implies O \sqsubseteq_\psi \mathbb{O}$

Proof: Theorems 1 - 4 are proved in Lemmas 4.19-4.16. Here we focus on the proof of TOPRULE. Suppose $O = (\eta_a, \eta_i, \theta)$ and $\mathbb{O} = (A, (\varphi, \varepsilon, \chi))$, then from the definition of $O \sqsubseteq_\psi \mathbb{O}$, we need to prove that:

for any $A, \Delta, \Sigma, \Delta, t, T$, if $\text{Match}(\psi, (T, \Delta, A, t), (T, \Delta, \Sigma))$ then

$$((A, O), W) \preceq ((A, \mathbb{O}), \mathbb{W})$$

then from Lemma 4.9, we need prove that

$$((A, O), W) \preceq ((A, \mathbb{O}), \mathbb{W})$$

from the definition of Match we know there exists T, Δ, t, A, Σ such that,

$$W = (T, \Delta, A, t_c) \quad \mathbb{W} = (T, \Delta, \Sigma) \quad (\psi \ A \ \Sigma)$$

suppose that $A = ((G, H, M), \text{isr}, \pi)$, then from Lemma 4.10, we need to prove that there exists $M_s, \Sigma_s, M_1, \Sigma_1, M_2, \Sigma_2, \dots, M_n, \Sigma_n, C_{lc}, C_{hc}$ such that:

- (1) $\Gamma; \chi; I \models \theta : \varepsilon \quad \chi; I \models \eta_i : \Gamma \quad \chi; I \models \eta_i : \Gamma$
- (2) $M = M_1 \uplus M_2 \uplus \dots \uplus M_n \uplus M_s, \Sigma = \Sigma_1 \uplus \Sigma_2 \uplus \dots \uplus \Sigma_n \uplus \Sigma_s$
- (3) $A|_{t_c} = \sigma_c, (\sigma_c \triangleleft M_s, \Sigma_s, _) \models [I]$
- (4) $P; \mathbb{P}; I; p \models (C_{lc}, \sigma_c \triangleleft M_c) \preceq (C_{hc}, \Sigma_c) \quad T_l(t_c) = C_{lc} \quad T_h(t_c) = C_{hc}$
- (5) for any $i, \sigma_i, \sigma_r, \Sigma_r, C_{li}, C_{hi}, i \neq c, T_l(t_i) = C_{li}, T_h(t_i) = C_{hi} \sigma_i = (A|_{t_i}) \triangleleft M_i, (\sigma_r, \Sigma_r, _) \models \text{SWINV}(I), \sigma_i \perp \sigma_r, \Sigma_i \perp \Sigma_r$, then $P; \mathbb{P}; I; p \models (C_{li}, \sigma_i \uplus \sigma_r) \preceq (C_{hi}, \Sigma_i \uplus \Sigma_r)$

where $((G, E, M), isr, \delta) \triangleleft M' \stackrel{\text{def}}{=} ((G, E, M'), isr, \delta)$. From TOPRULE and Lemma 4.19 - 4.16, we can trivially know that (1) holds, then from TOPRULE we know that

$$[\psi] \Rightarrow I[0, N] * \text{OS}[\bar{0}, 1, \text{nil}, \text{nil}] * \text{empE} \quad (1)$$

then from the semantics of the assertion we know that exist

$$M_1 = M_2 = \dots = M_n = \emptyset, \quad \Sigma_1 = \Sigma_2 = \dots = \Sigma_n = \emptyset$$

and

$$M_s = M, \quad \Sigma_s = \Sigma$$

such that (2) and (3) hold.

Then from the definition of **Match**, we know that

$$\forall i, \exists s_i, C_{li} = C_{hi} = (s_i, (\circ, \bullet))$$

Then from Lemma 4.11 we know that (4) holds. From the definition **SWINV**(I) and (I) we know that, $((A|_{t_i}) \triangleleft \emptyset, \emptyset, -) \models \text{SWINV}(I)$ and we know that $A|_{t_i} = (A|_{t_i}) \uplus ((A|_{t_i}) \triangleleft \emptyset)$, then for (5), we only need to prove that for any i , if $i \neq c$, then

$$P; \mathbb{P}; I; p \models (C_{li}, A|_{t_i}) \preceq (C_{hi}, \Sigma_i)$$

and it can be easily proved by Lemma 4.11.

Lemma 4.13. For all $\Gamma, \chi, I, r, p_i, p, s$ and $q, \Gamma; \chi; I; r; p_i \vdash \{p\} s \{q\} \Rightarrow \Gamma; \chi; I; r; p_i \models \{p\} s \{q\}$

Proof: First induction over the inference rules, for the compositional cases, like while, we prove them by co-induction. The proof details can be found in [?].

Lemma 4.14. For all Γ, χ, I, θ and $\varepsilon, \Gamma; \chi; I \vdash \theta : \varepsilon \Rightarrow \Gamma; \chi; I \models \theta : \varepsilon$

Proof: By Def. 4.4 and Lemma 4.19.

Lemma 4.15. For all Γ, χ, I, η_a and $\varphi, \Gamma; \chi; I \vdash \eta_a : \varphi \Rightarrow \Gamma; \chi; I \models \eta_a : \varphi$

Proof: By Def. 4.5 and Lemma 4.19.

Lemma 4.16. For all Γ, χ, I and $\eta, \chi; I \vdash \eta : \Gamma \Rightarrow \chi; I \models \eta : \Gamma$

Proof: By Def. 4.6 and Lemma 4.19.

Definition 4.17 (Simulation).

$$\Gamma; \chi; I; r; p_i; r' \models (C, \sigma) \preceq (\mathfrak{s}, \Sigma)$$

holds, whenever:

- **Normal Steps:** for any $\eta, C', M_s, M_f, \Sigma_s, \sigma_1$ and σ'_1 , if $\neg \text{lsFcall}(C), (\sigma|_{M_s}, \Sigma_s, \mathfrak{s}) \models [I], \sigma_1 = \sigma \uplus_\triangleleft M_s \uplus_\triangleleft M_f, \Sigma \perp \Sigma_s$ and $\eta \vdash (C, \sigma_1) \bullet \text{L} \rightarrow \triangleright (C', \sigma'_1)$, then there exist $\Sigma'_s, \mathfrak{s}', \Sigma', \sigma'$ and M'_s , such that the followings hold:

- $\sigma'_1 = \sigma' \uplus_{\triangleleft} M'_s \uplus_{\triangleleft} M_f, (\sigma|_{M'_s}, \Sigma'_s, \mathfrak{s}') \models [I],$
 - $(\mathfrak{s}, \Sigma \uplus \Sigma_s) \bullet\text{--}H\text{--}\overset{*}{\triangleright} (\mathfrak{s}', \Sigma' \uplus \Sigma'_s),$
 - $\Gamma; \chi; I; r; p_i; r' \models (C', \sigma') \preceq (\mathfrak{s}', \Sigma').$
- **Function Call:** for any $M_s, \Sigma_s, \kappa_s, f, \bar{v}$ and \mathcal{T} , if $C = (\mathbf{fexec}(f, \bar{v}), (\circ, \kappa_s)), (\sigma.m.M) \perp M_s, (\sigma|_{M_s}, \Sigma_s, \mathfrak{s}) \models [I]$ and $\Sigma \perp \Sigma_s$, then there exist $\sigma_1, M_f, \Sigma_1, \Sigma_f, \Sigma', \Sigma'_s, \mathfrak{s}', fp, fq, \tau$ and \mathcal{L} , such that the followings hold:
- $\Gamma(f) = (fp, fq, (\tau, \mathcal{T})),$
 - $(\mathfrak{s}, \Sigma \uplus \Sigma_s) \bullet\text{--}H\text{--}\overset{*}{\triangleright} (\mathfrak{s}', \Sigma' \uplus \Sigma'_s), (\sigma|_{M_s}, \Sigma'_s, \mathfrak{s}') \models [I],$
 - $\sigma = \sigma_1 \uplus_{\triangleleft} M_f, \Sigma' = \Sigma_1 \uplus \Sigma_f, (\sigma_1, \Sigma_1, \mathfrak{s}') \models fp(\text{rev}(\bar{v})) \mathcal{L},$
 - for any $\sigma', \sigma'_1, \Sigma'', \Sigma'_1, \hat{v}$ and \mathfrak{s}'' , if $\sigma.m.G = \sigma'.m.G, \sigma.m.E = \sigma'.m.E,$
 $(\sigma'_1, \Sigma'_1, \mathfrak{s}'') \models fq(\text{rev}(\bar{v})) \hat{v} \mathcal{L},$
 $\sigma' = \sigma'_1 \uplus_{\triangleleft} M_f$, and $\Sigma'' = \Sigma'_1 \uplus \Sigma_f$, then
 $\Gamma; \chi; I; r; p_i; r' \models ((\mathbf{skip} \hat{v}, (\circ, \kappa_s)), \sigma') \preceq (\mathfrak{s}'', \Sigma'').$
- **Context Switch:** for any M_s, κ_s, x and Σ_s , if $(\sigma.m.M) \perp M_s, C = (\mathbf{switch} x, (\circ, \kappa_s)), (\sigma|_{M_s}, \Sigma_s, \mathfrak{s}) \models [I]$ and $\Sigma \perp \Sigma_s$, then there exist $\sigma_1, M, \Sigma', \Sigma'_s$ and \mathfrak{s}' , such that the followings hold:
- $(\mathfrak{s}, \Sigma \uplus \Sigma_s) \bullet\text{--}H\text{--}\overset{*}{\triangleright} (\mathbf{sched}; \mathfrak{s}', \Sigma' \uplus \Sigma'_s),$
 - $(\sigma|_{M_s}, \Sigma'_s, \mathfrak{s}') \models [I], \sigma = \sigma_1 \uplus_{\triangleleft} M,$
 - $(\sigma|_M, \Sigma', \mathfrak{s}') \models \text{SWINV}(I) \wedge \text{SWPRE}(\chi, x)$
 - for any M', σ' and Σ'' , if $\sigma' = \sigma_1 \uplus_{\triangleleft} M'$ and $(\sigma'|_{M'}, \Sigma'', \mathfrak{s}') \models \text{SWINV}(I),$
then $\Gamma; \chi; I; r; p_i; r' \models ((\mathbf{skip} \perp, (\circ, \kappa_s)), \sigma') \preceq (\mathfrak{s}', \Sigma'').$
- **Skip:** for any \hat{v}, M_s and Σ_s , if $C = (\mathbf{skip} \hat{v}, (\circ, \bullet)), (\sigma.m.M) \perp M_s, (\sigma|_{M_s}, \Sigma_s, \mathfrak{s}) \models [I]$ and $\Sigma \perp \Sigma_s$, then there exist Σ', Σ'_s and \mathfrak{s}' , such that $(\mathfrak{s}, \Sigma \uplus \Sigma_s) \bullet\text{--}H\text{--}\overset{*}{\triangleright} (\mathfrak{s}', \Sigma' \uplus \Sigma'_s),$
 $(\sigma|_{M_s}, \Sigma'_s, \mathfrak{s}') \models [I],$ and $(\sigma, \Sigma', \mathfrak{s}') \models r' \hat{v}.$
- **Return:** for any κ_s, M_s and Σ_s , if $C = (\mathbf{return}, (\circ, \kappa_s)), \lfloor \kappa_s \rfloor = \perp, \lfloor \kappa_s \rfloor_c = \perp,$
 $(\sigma.m.M) \perp M_s, (\sigma|_{M_s}, \Sigma_s, \mathfrak{s}) \models [I]$ and $\Sigma \perp \Sigma_s$, then there exist Σ', Σ'_s
and \mathfrak{s}' , such that $(\mathfrak{s}, \Sigma \uplus \Sigma_s) \bullet\text{--}H\text{--}\overset{*}{\triangleright} (\mathfrak{s}', \Sigma' \uplus \Sigma'_s), (\sigma|_{M_s}, \Sigma'_s, \mathfrak{s}') \models [I],$ and
 $(\sigma, \Sigma', \mathfrak{s}') \models r \perp.$
- **ReturnE:** for any v, κ_s, M_s and Σ_s , if $\lfloor \kappa_s \rfloor = \perp, \lfloor \kappa_s \rfloor_c = \perp, C = (\mathbf{skip} v, (\circ, (\mathbf{return} _). \kappa_s)), (\sigma.m.M) \perp M_s,$
 $(\sigma|_{M_s}, \Sigma_s, \mathfrak{s}) \models [I]$ and $\Sigma \perp \Sigma_s$, then there exist Σ', Σ'_s and \mathfrak{s}' , such that
 $(\mathfrak{s}, \Sigma \uplus \Sigma_s) \bullet\text{--}H\text{--}\overset{*}{\triangleright} (\mathfrak{s}', \Sigma' \uplus \Sigma'_s),$
 $(\sigma|_{M_s}, \Sigma'_s, \mathfrak{s}') \models [I],$ and $(\sigma, \Sigma', \mathfrak{s}') \models r v.$
- **IRet:** for any κ_s, M_s and Σ_s , if $C = (\mathbf{ixext}, (\circ, \kappa_s)), (\sigma.m.M) \perp M_s, (\sigma|_{M_s}, \Sigma_s, \mathfrak{s}) \models [I]$ and $\Sigma \perp \Sigma_s$, then there exist Σ', Σ'_s and \mathfrak{s}' , such that $(\mathfrak{s}, \Sigma \uplus \Sigma_s) \bullet\text{--}H\text{--}\overset{*}{\triangleright} (\mathfrak{s}', \Sigma' \uplus \Sigma'_s),$
 $(\sigma|_{M_s}, \Sigma'_s, \mathfrak{s}') \models [I],$ and $(\sigma, \Sigma', \mathfrak{s}') \models q \hat{v}.$
- **Abort:** for any η, Σ_s, M_s, M_f and σ' , if $\neg \text{IsFcall}(C), \sigma' = \sigma \uplus_{\triangleleft} M_s \uplus_{\triangleleft} M_f,$
 $(\sigma|_{M_s}, \Sigma_s, \mathfrak{s}) \models [I]$ and $\Sigma \perp \Sigma_s$, then
 $\neg(\eta \vdash (C, \sigma') \bullet\text{--}L\text{--}\triangleright \mathbf{abort})$

When executing the client codes we require that each step of low-level and high-level will keep the client state and code the same. $\mathbf{fexec}(f, \bar{v}\mathcal{T})$ is the entry point of the OS kernel(f is a OS API).

For each step of the kernel, we require it to be safe, and corresponds to some steps of high level code, as shown in **Normal Steps** in Def. 4.17. Since the interrupts may come at any time(if it is enable), we will get a new well-formed shared resource which is specified by $[I]$ at each step.

To support modular reasoning of internal function calls, at the entry point of a internal function, we use the function specification stored in Γ to avoid step into the function body, as shown in the **Function Steps** in Def. 4.17.

For the **Context Switch** case, we require that the isr is clear, ie is 0, low-level and high-level can be scheduled to the same task, and all the shared resources from I_0 to I_n are well-formed. Because the interrupt is disable, so there exists some shared resources occupied by the current task, we specify them with $SWINV(I)$ to make sure that it is well-formed.

Definition 4.18 (Judgment Semantics). $\Gamma; \chi; I; r; p_i \models \{p\}s\{q\}$

holds, iff for any σ , Σ and \mathbb{s} , if $(\sigma, \Sigma, \mathbb{s}) \models p$, then

$\Gamma; \chi; I; r; p_i; \lambda_{-}q \models ((s, (\circ, \bullet)), \sigma) \preceq (\mathbb{s}, \Sigma)$

Theorem 4.19 (Soundness of Inference Rules).

$\Gamma; \chi; I; r; p_i \vdash \{p\}s\{q\} \implies \Gamma; \chi; I; r; p_i \models \{p\}s\{q\}$

Definition 4.20 (Well-Formed Interrupts). $\Gamma; \chi; I \vdash \theta : \varepsilon$ holds, iff for any k , isr , is , G , p and p_i , if $p = \text{BldltpPre}(k, \varepsilon, isr, is, I)$, $p_i = \text{BldltpRet}(k, isr, is, I)$, then there exists s , such that $\theta(k) = s$ and $\Gamma; \chi; I; \lambda_{-}\text{false}; p_i \vdash \{p\}s\{\text{false}\}$

Definition 4.21 (Well-Formed APIs). $\Gamma; \chi; I \vdash \eta_a : \varphi$ holds, iff for any f , \bar{v} , p and r , if $f \in \text{dom}(\varphi)$, $p = \text{BuildAPIPre}(\eta_a, f, \varphi, \bar{v})$, $r = \text{BuildAPIRet}(\eta_a, f)\bar{v}$, then there exists s , such that $\eta_a(f) = (-, -, -, s)$ and $\Gamma; \chi; I; r; \text{false} \vdash \{p\}s\{\text{false}\}$

Definition 4.22 (Well-Formed Internal Functions). $\chi; I \vdash \eta_i : \Gamma$

holds, iff $\text{dom}(\Gamma) = \text{dom}(\eta_i)$ and for any f , fp and fq if $\Gamma(f) = (fp, fq)$, then there exist τ , \mathcal{D}_1 , \mathcal{D}_2 and s , such that

- $\eta_i(f) = (\tau, \mathcal{D}_1, \mathcal{D}_2, s)$, and
- for any \bar{v} , p , r , and \mathcal{L} , if $p = \text{BuildFunPre}(\eta_i, f, \bar{v}, fp)$ and $r = \text{BuildFunRet}(\eta_i, f, \bar{v})fq$, then $\Gamma; \chi; I; r; \text{false} \vdash \{p\}s\{\text{false}\}$

Theorem 4.23 (Verifying OS Correctness). For any O , \mathbb{O} and ψ , if there exist Γ , I , $\chi, \eta_a, \eta_i, \theta$, φ, ε and χ , such that $O = (\eta_a, \eta_i, \theta)$, $\mathbb{O} = (\varphi, \varepsilon, \chi)$ and the followings hold:

- **Verifying Interrupts:** $\Gamma; \chi; I \vdash \theta : \varepsilon$,
- **Verifying APIs:** $\Gamma; \chi; I \vdash \eta_a : \varphi$,
- **Verifying Internal Functions:** $\chi; I \vdash \eta_i : \Gamma$,
- **Side Conditions :** $\text{Good}(I, \chi)$ and for any Λ and Σ , if $\psi \wedge \Sigma$, then the followings hold :
 - $\forall t. (\Lambda|_t, \Sigma, -) \models [I] * \text{OS}[\text{empisr}, 1, \text{nil}, \text{nil}] * \text{LVar}(\text{nil})$,

then $O \sqsubseteq_{\psi} \mathbb{O}$.

```

L1  OS_ENTER_CRITICAL();
L2  ticks = OStime;
L3  OS_EXIT_CRITICAL();
L4  return (ticks);

Lemma OStimeGetRight: forall r p,
  Some r = BuildRetA' api TimeGet tmgetspec nil ->
  Some p = BuildPreA' api TimeGet tmgetspec nil ->
    exists t d1 d2 s, api TimeGet = Some (t, d1, d2, s) /\
      {|F, X, I, r, Afalse|} |- {|p|} s {|Afalse|}.
Proof.
  init spec.          (*Initialize Specification*)
  hoare forward.       (*Forward L1*)
  hoare unfold pre.    (*Prepare the precondition*)
  hoare forward.       (*Forward L2 *)
  hoare absconsq.      (*The high-level step*)
  eapply OStimeGet_high_level_step; eauto.
  hoare forward.       (*Forward L3*)
  unfold AOStime.
  sep auto.           (*Solve side conditions*)
  eauto.
  hoare forward.       (*Forward L4*)
Qed.

```

Fig. 26. Verifying OStimeGet() with Tactics

4.5 Coq Tactics

We develop a set of practical tactics in Coq based on the rules of the refinement logic, including “sep auto” for proving relational assertions, “hoare forward” for proving refinement judgments, and some domain-specific tactics for proving the properties of integers and bitmaps. Here we omit the implementation details which can be seen in [28].

To demonstrate the efficiency of our tactics, we make two versions of proofs [28] for the API (OStimeGet), which has only 4 lines code and uses a critical region to read the global clock and return its value. The one using our tactics only needs 11 lines of proof scripts, while the other one using the primitive tactics provided by Coq needs more than 400 lines.

Another advantage of our tactics is that they can help us to smoothly reason about the spatial parts and extract lemmas that are independent of program contexts for verifying functionality of kernels. Users with little knowledge about our framework can prove the code with our tactics.

Because all the APIs are proved in the similar procedure as shown in Fig. 26, first we initialize specifications with the tactics “init spec”, and we obtain the Hoare triple with pre- and post- conditions instantiated, next we apply the “hoare forward” tactic step by step till the point before exiting the critical region. Then the high-level specification code must be executed to reestablish the

invariant over the related states and exit. Before that we may need to derive the necessary premises for safely executing the specification code in terms of the states satisfying invariants when entering the critical section. By applying the tactic “`hoare absq`”, it allows us to change the precondition by applying a lemma of proving the safe execution of the current specification code with one-or-multiple steps. Note that the initial specification code usually contains many statements of non-deterministic choice, and it is users’ responsibility to choose a correct branch for accomplishing the proofs. The tactic “`sep auto`” is used to solve the generated side conditions about entailment of relational assertions.

Then the major proof efforts lie in proving the extracted lemmas, which are related to the functional correctness of kernel APIs. For instance, $\mu\text{C}/\text{OS-II}$ supports 64 tasks with priorities from 0 to 63, and the scheduling algorithm calculates the highest priority of ready tasks using bit operations over the ready table, and we need to prove the lemma like “ $0 \leq x < 64 \rightarrow 0 \leq (x \gg 3) < 8$ ”, which involves bit operations over integers limited in small finite domains, such as 0-8, 0-64 and 0-256. Most of the mathematical properties required by $\mu\text{C}/\text{OS-II}$ are limited with small finite domains like this. We develop a domain-specific tactic called “`mauto`” to automatically prove these mathematical properties. “`mauto`” is implemented by brute force iterating all the possible inputs and solving each subgoal with “`omega`”.

5 Proving Priority-Inversion-Freedom

Definition of priority-inversion-freedom (PIF). PIF is an important property in real-time systems, but there is no standard formal definition for it. Although there have been efforts trying to formalize and verify it, their definitions all have serious problems. For instance, Def. 5.1 is formal definition of priority inversion given in earlier work [6]. It says priority inversion occurs if there is a higher priority task t waiting directly or indirectly for a lower priority task t' .

Definition 5.1 (Priority Inversion). $\text{PI}(\Sigma)$ holds, iff there exist t and t' such that $t \xrightarrow{\Sigma}^+ t'$ and $\text{CurPr}(t', \Sigma) < \text{CurPr}(t, \Sigma)$.

Here the waiting chain $t \xrightarrow{\Sigma}^+ t'$ is the transitive closure of the waiting relation $t \xrightarrow{\Sigma} t'$, saying t waits for the resource owned by t' . $\text{CurPr}(t, \Sigma)$ returns the current priority of t . Then PIF can be simply defined as the negation of $\text{PI}(\Sigma)$ over any state Σ on the execution sequence.

Although it looks simple and intuitive, it cannot be applied to classic real-world algorithms for PIF, *e.g.*, priority ceiling and priority inheritance [25], because they need to dynamically change the priority of tasks. Since the definition refers to the current priority of tasks, its meaning crucially depends on the algorithms, which becomes difficult to understand.

A reasonable definition must be based on the original priorities assigned by the programmer, reflecting the actual degree of urgency. Below we give a new definition for PIF.

Definition 5.2 (Priority Inversion Freedom). $\text{PIF}(\Sigma)$ holds, iff for any t , t_c , pr and pr_c , if $t \neq t_c$, $t_c = \text{CurTask}(\Sigma)$, $pr = \text{OrgPr}(t, \Sigma)$, $pr_c = \text{OrgPr}(t_c, \Sigma)$, $\text{IsWait}(t, \Sigma)$ and $\neg \text{IsOwner}(t_c, \Sigma)$, then $pr \preceq pr_c$.

It says, if the current task t_c does not own any shared resources, then its *original* priority should be higher than (or equal to) any other waiting tasks t . Here $\text{OrgPr}(t, \Sigma)$ represents t 's original priority configured by users. $\text{IsWait}(t, \Sigma)$ is defined as $\exists t'. t \xrightarrow{\Sigma}^+ t'$, and $\neg \text{IsOwner}(t_c, \Sigma)$ means that the task t_c does not own any shared resources (*e.g.*, mutexes).

The definition essentially says that the current task can have a lower priority than waiting ones only if it holds resources that might be needed by others (so we want to run it first to make the resource available as soon as possible). Note that if each task eventually releases its shared resource (*i.e.*, there is no deadlock), the waiting task with higher priority will be eventually released and executed. Therefore the definition prevents unbounded priority inversion [25].

Here we give an informal argument as a justification of our new definition. Suppose we have tasks $A_1, \dots, A_m, B_1, \dots, B_n$ and C , whose original priorities have the order: $A_1 > \dots > A_m > B_1 > \dots > B_n > C$. B_i is the current running task, and A_m is waiting for the resource acquired by C . According to our PIF definition, we are able to show that A_m will not be permanently delayed by the lower-priority tasks $(B_1, B_2, \dots, B_n, C)$. That is, our PIF definition indeed prevents unbounded priority inversion. Let us consider the following two cases for the current running task B_i :

- (1) If B_i does not own any resource, then the scenario obviously violates our PIF definition, which requires that the current running task should have higher original priority than the waiting one A_m ;
- (2) If B_i owns some resources, with the assumption of termination of critical regions, B_i eventually exits its outmost critical region, then we have the following three different cases:
 - If B_i is still the current running task, then it violates PIF as in case (1);
 - If we switch to B_j , then one of the following two cases holds:
 - B_j does not own resources, and it violates PIF as in (1);
 - B_j owns some shared resources. Then we get to the same scenario as (2) above where B_j now plays the role of B_i , but now the number of tasks that are inside critical regions is decreased by 1. By induction over the number of tasks inside of critical regions, A_m can only be delayed by the lower-priority tasks $(B_1, B_2, \dots, B_n, C)$ with bounded periods (bounded number of steps of running critical regions).
 - If we switch to A_k instead, then it is OK because A_k has higher original priority than A_m .

We demonstrate the difference between the two notions by showing in Fig. 27 an example violating PIF. It is constructed based on the mutex implementation in $\mu\text{C}/\text{OS-II}$, which is implemented with a simplified priority ceiling protocol [25] and cannot prevent priority inversions when there is nested use of mutexes. The counterexample justifies this limitation of nested mutex PIF failure of $\mu\text{C}/\text{OS-II}$.

Before going through the example, we first explain the protocol used for $\mu C/OS-II$ mutex. It provides two APIs to acquire ($OSMutexPend(S)$) and release ($OSMutexPost(S)$) the mutex, $P(S)$ and $V(S)$ for short. Each mutex S has a unique priority, and it saves its owner task's identifier and current priority when S is acquired. When a task t executes $P(S)$, it executes the following steps: (1) t 's current priority must be lower than that of S ; (2) if S is available, t successfully acquires S , sets S 's owner to t , and saves t 's current priority in S for the future recovery; and (3) t is blocked if S is already owned by another task t' . If t 's current priority is higher than that of t' , we lift the current priority of t' to S 's priority. We omit the behavior of $V(S)$ here, which is not needed to understand the counterexample.

In Fig. 27, we show four tasks (A, B, C and D) and two mutexes (S_1 and S_2). They all have their original priorities configured by users, following the order of $Pr_D < Pr_C < Pr_{S_2} < Pr_B < Pr_A < Pr_{S_1}$. Suppose A, B and C are blocked and waiting for the timer interrupt to wake up, and D is the only running task. First D acquires S_2 , then C wakes up and preempts D . C acquires S_1 , and then attempts to acquire S_2 owned by D , thus gets blocked. Since $Pr_D < Pr_C$, following the aforementioned protocol we lift D 's priority to that of S_2 (Pr_{S_2}). Then A wakes up and becomes the highest priority task to run. It tries to acquire S_1 (now owned by C) and gets blocked. Then we lift C 's priority to that of S_1 (Pr_{S_1}). Finally B becomes ready and has the highest priority to run. In each step we also show the pairs of the current priority and the task status for corresponding tasks. Now we have:

- **Violation of the old definition:** C waits for D , but C 's current priority Pr_{S_1} is higher than that of D (Pr_{S_2}), shown in the dashed box in Fig.27;
- **Violation of our new definition:** B does not own any mutexes, and the waiting task A 's original priority Pr_A is higher than that of B (Pr_B), shown in the solid box.

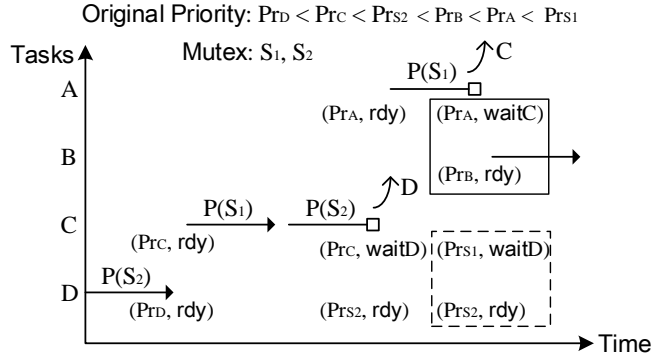


Fig. 27. Violating PIF with Mutex in $\mu C/OS-II$

In this example, unlike C and D that are going to release mutexes in bounded number of steps, B does not own any mutexes and may run forever, therefore the most urgent task A may never get a chance to run. This situation is called unbounded priority inversion [25] and should be forbidden in real-time systems. On the other hand, if we only have A, C and D , A only needs to wait for D

$$\begin{aligned}
\text{CurTask}(\Sigma) &\stackrel{\text{def}}{=} \Sigma(\text{tcbls})(\text{ctid}) & \text{CurPr}(t, \Sigma) &\stackrel{\text{def}}{=} \begin{cases} pr & \text{if } \Sigma(\text{tcbls})(t) = (pr, -) \\ \perp & \text{otherwise} \end{cases} \\
\text{Own}(t, \text{eid}, \Sigma) &\stackrel{\text{def}}{=} t \in \text{dom}(\Sigma(\text{tcbls})) \wedge \Sigma(\text{ecbls})(\text{eid}) = (-, (t, -)) \\
\text{IsOwner}(t, \Sigma) &\stackrel{\text{def}}{=} \exists \text{eid}. \text{Own}(t, \text{eid}, \Sigma) & \text{OrgPr}(t, \Sigma) &\stackrel{\text{def}}{=} \begin{cases} pr & \text{if } \exists \text{eid}. \text{Own}(t, \text{eid}, \Sigma) \wedge \\ & \Sigma(\text{ecbls})(\text{eid}) = (-, (-, pr)) \\ pr & \text{if } (\neg \text{IsOwner}(t, \Sigma)) \wedge \\ & \Sigma(\text{tcbls})(t) = (pr, -) \\ \perp & \text{otherwise} \end{cases} \\
\text{Wait}(t, \text{eid}, \Sigma) &\stackrel{\text{def}}{=} \Sigma(\text{tcbls})(t) = (-, \text{wait}(\text{mtx}(\text{eid}), -)) \\
t \xrightarrow{\Sigma} t' &\stackrel{\text{def}}{=} \exists \text{eid}. \text{Wait}(t, \text{eid}, \Sigma) \wedge \text{Own}(t', \text{eid}, \Sigma) \wedge t \neq t' \\
\text{IsWait}(t, \Sigma) &\stackrel{\text{def}}{=} \exists t'. t \xrightarrow{\Sigma} t' & \text{IsRdy}(t, \Sigma) &\stackrel{\text{def}}{=} \Sigma(\text{tcbls})(t) = (-, \text{rdy})
\end{aligned}$$

Fig. 28. Auxiliary Definitions for PIF

and C to release their mutexes before it is unblocked. There is no unbounded priority inversion. This case is allowed in our new definition, but disallowed in the old one (as shown in the dashed box).

The above counterexample happens because the priority lifting implemented in $\mu\text{C}/\text{OS-II}$ fails to consider nested use of mutexes. Below we show it could satisfy PIF if we disallow nested use of mutexes.

Proving PIF of Mutex in $\mu\text{C}/\text{OS-II}$. We prove that the mutex in $\mu\text{C}/\text{OS-II}$ ensures both the old and our new PIF definitions if there are no nested use of mutexes. Some auxiliary definitions are given in Fig.28 based on the mutex implementation of $\mu\text{C}/\text{OS-II}$. Here we present Theorem 5.3 based on our PIF definition.

Theorem 5.3 (PIF without Nested Use of Mutexes). If $\text{Init}(\Sigma), (A, \mathbb{O}_{\mu\text{C}/\text{OS-II}}) \vdash (T, \Delta, \Sigma) =_{H \Rightarrow^*} (T', \Delta', \Sigma')$, $\text{NoNCR}(A, \Sigma, T, \Delta)$, and $\text{SchedProp}(\Sigma')$, then $\text{PIF}(\Sigma')$. Where $\text{Init}(\Sigma) \stackrel{\text{def}}{=} \forall t \in \text{dom}(\Sigma(\text{tcbls})). \neg \text{IsWait}(t, \Sigma)$.

It says, for any application code A , client state Δ and kernel abstract state Σ , if initially there are no tasks waiting for mutexes ($\text{Init}(\Sigma)$), there is no nested use of mutexes in the program ($\text{NoNCR}(A, \Sigma, T, \Delta)$), then for any T' , Δ' and Σ' generated during the execution, if Σ' is consistent with the priority-based scheduling ($\text{SchedProp}(\Sigma')$), then it must satisfy PIF.

The following definition $\text{NoNCR}(\eta_c, \Sigma, T, \Delta)$ is semantically defined by excluding program states that the execution of nested use of mutexes might get to.

Definition 5.4 (No Nested). $\text{NoNCR}(\eta_c, \Sigma, T, \Delta)$ holds, iff for any Σ' , T' and Δ' , if $(A, \mathbb{O}_{\mu\text{C}/\text{OS-II}}) \vdash (T, \Delta, \Sigma) =_{H \Rightarrow^*} (T', \Delta', \Sigma')$, then $\text{NNest}(\Sigma')$ defined as below.

$$\text{NNest}(\Sigma) \stackrel{\text{def}}{=} \forall t, eid. (t \in \text{dom}(\Sigma(\text{tcbls})) \wedge \text{Own}(t, eid, \Sigma)) \rightarrow \\ \neg(\text{IsWait}(t, \Sigma) \vee \exists eid'. eid' \neq eid \wedge \text{Own}(t, eid', \Sigma))$$

$\text{SchedProp}(\Sigma')$ given in Def.5.5 requires that the current running task always has the highest priority among all the ready tasks. It can be guaranteed by the highest-priority-based scheduling strategy of $\mu\text{C}/\text{OS-II}$. Here we use a simplified $\mathbb{O}_{\mu\text{C}/\text{OS-II}}$ that contains the PIF mutex as the only APIs. $\text{IsRdy}(t, \Sigma)$ means that the task t is ready in Σ . The proof is formalized in Coq.

Definition 5.5 (Highest-Priority-Based Scheduling).

$\text{SchedProp}(\Sigma)$ holds, iff for any t_c and pr_c , if $t_c = \text{CurTask}(\Sigma)$ and $pr_c = \text{CurPr}(t_c, \Sigma)$, then $\text{IsRdy}(t_c, \Sigma)$ and for any t and pr , if $t \neq t_c$, $pr = \text{CurPr}(t, \Sigma)$ and $\text{IsRdy}(t, \Sigma)$, then $pr \preceq pr_c$.

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

We have applied our framework to verify key modules (around 1300 lines of C code without counting comments and empty lines) of $\mu\text{C}/\text{OS-II}$ V2.52, including the scheduler, the timer interrupt handler, mutexes, message queues, mail boxes, semaphores, and the time management. These 1300 lines of C code verified in our framework correspond to around 3250 lines of code in their original format (with comments and empty lines) in the source files of $\mu\text{C}/\text{OS-II}$, including “ucos_ii.h”, “os.q.c”, “os.sem.c”, “os.mbox.c”, “os.mutex.c”, “os.time.c”, “os.core.c” and “os.cpu.a.c”. The verified modules cover 63% of the frequently used APIs and internal functions [2]. We ignore some synchronization APIs which have similar functionality as the verified ones. Verification of task creation/deletion is still ongoing work based on the presented framework.

Invariants. Because we are trying to verify a third-party developed OS kernel, we need to spell out the invariants based on our understandings to kernel implementations. It is more difficult than specifying a newly developed kernel with verification in mind.

For instance, the most important data structure in $\mu\text{C}/\text{OS-II}$ is a double-linked-list of task control blocks (TCB). It is shared among the interrupt handlers and tasks. The following assertion is used to specify the TCB list, and we use it to demonstrate the efforts to define the invariants for $\mu\text{C}/\text{OS-II}$.

$$\text{tcddlseg}(t, t_1, t_2, t_3, \mathcal{L}) * \text{tcbls} \mapsto \alpha * \text{OSRdyTbl}(\bar{v}) * R_{\text{tcbls}}(t, \mathcal{L}, \alpha, \bar{v})$$

$\text{tcddlseg}(t, t_1, t_2, t_3, \mathcal{L})$ is inductively defined over \mathcal{L} to specify a segment of double-linked-list of TCBs. \mathcal{L} is a list of \bar{v} containing all the values stored in the TCB list segment, and t is the head pointer (the first task’s id). The abstract TCB list α maps task identifiers to their abstract TCBs. In addition to the spatial parts, we also need to define the relation between the low-level TCB list \mathcal{L} and the high-level abstract TCB list α , denoted as $R_{\text{tcbls}}(t, \mathcal{L}, \alpha, \bar{v})$. The relation is parametrized with a value list \bar{v} containing the values stored in the ready table OSRdyTbl , which is a bitmap for recording the status of tasks. Thus the status of tasks in abstract TCBs depends on both \mathcal{L} and \bar{v} . There are two different ways

Framework	Coq lines	Verified Modules	lines of C	Coq lines
Basic Libraries	32061	Global Declarations	187	-
Machine & Logic	23095	Message Queue	240	4537
Automated Tactics	21050	Semaphore	166	2441
Total	76206	Mailbox	171	3326
Certified $\mu C/OS-II$	Coq lines	Mutex	301	17331
C Code Definitions	1824	Time Management	39	861
Specifications	6012	Timer Interrupt	17	443
Priority Inversion Freedom	9570	Internal Functions	195	5447
Libraries for $\mu C/OS-II$	62085	Final Theorems	-	501
Auto. Generated Code	25357	Total	1316	34887
Total	104848			

Table 1. The Verification Package

of defining the relation $R_{tcbls}(t, \mathcal{L}, \alpha, \bar{v})$, one is to inductively define as below and the other way is to define it with quantifications.

$$R_{tcbls}(t, \bar{v}' :: \mathcal{L}, \{t \rightsquigarrow \mathbf{abstcb}\} \uplus \alpha, \bar{v}) \stackrel{\text{def}}{=} RL_{tcb}(t, \bar{v}', \mathbf{abstcb}, \bar{v}) \wedge R_{tcbls}(\text{next}(\bar{v}'), \mathcal{L}, \alpha, \bar{v})$$

We find that the inductive definition can greatly simply our proofs. The automated technique [23] for separation logic also shows that inductively defining the spatial parts and data constraints over the same inductive data lets automated proofs for some complex data structures become possible.

Here $RL_{tcb}(t, \bar{v}', \mathbf{abstcb}, \bar{v})$ describes the correspondence relation between the low-level TCB and its high-level abstract TCB. Note that `OSRdyTbl`'s value list \bar{v} cannot be partitioned into each TCB with the inductive definition. Then we need to prove that the update to `OSRdyTbl` for one task will not break the relation for the others. If we verify a kernel developed by ourself, we may simply avoid this burden by using a while-loop to search the TCB list for scheduling but not using the fast bitmap. Thus the gap between the low-level data structures and high-level abstract representation makes verification more difficult.

Modifications to the original code. Our verification is based on the original code with some minor modifications. For instance, the API `OSQPend(S)` is used to receive a message from a queue, and its original code does not check if the input pointer S points to a valid event control block, because it assumes that the client code always gets S by calling `OSQCreate()` (thus S should already be valid). We drop this assumption about the client code. Correspondingly we insert code that checks whether S is a valid pointer. If S is invalid a new error code is returned. Similar modifications are made to some other modules too. The reason for doing above modifications is that the contextual refinement proved in our verification framework assumes arbitrary client code, while kernels are usually implemented with assumptions over client code for efficiency.

Proof efforts. The Coq implementation consists of around 216,000 lines of code and proofs in Coq8.4pl6. Table 1 gives a break down of the number of lines for various components. Compiling the entire Coq package takes around 16 hours

on a machine with 3.6GHz cpu and 32G memory. The work takes us around 5.5 person years in total, including 4 person years for the framework and 1 person year for verifying the first $\mu\text{C}/\text{OS-II}$ module (Message Queue). With the facilities (tactics, libraries and invariants *etc.*) being stabilized, verifying the remaining modules (around 900 lines of C code) only takes us around 6 person months.

The most challenging part is to verify the timer interrupt handler, which traverses the entire TCB list and updates task status in each TCB block. It needs to access all the shared data structures in $\mu\text{C}/\text{OS-II}$. Several different updates to shared data structures make the loop invariant quite complicated.

Also verifying an existing OS kernel is more difficult than verifying a new one written for verification purpose. When verifying $\mu\text{C}/\text{OS-II}$ the major difficulty comes from the gap between the low-level concrete data structure and the high-level abstract representation. For instance, $\mu\text{C}/\text{OS-II}$ uses a smart bitmap algorithm to record whether a task is in the waiting queue. The implementation requires us to establish a subtle consistency relation between the low-level bitmap and the high-level abstract waiting queue. The verification would have been much simpler if the waiting queue is simply implemented as a linked list.

Coq tactics. Proof automation is essential to improve the productivity. We develop tactics for automatically proving relational separation logic assertions and generating verification conditions based on existing techniques [7, 21, 5]. They do forward reasoning for statements, including function calls and primitives entering and exiting critical regions, *etc.* Also some domain-specific tactics are implemented for individual data structures used in $\mu\text{C}/\text{OS-II}$, including ones for the arithmetic properties of *Int32* and bitmaps. Thanks to these tactics, the ratio of Coq proof scripts to the verified C code is around 26:1. Another advantage of the tactics is that they can extract lemmas independent of program contexts for verifying functionality of code. Users can verify code using the tactics without knowing much about the underlying framework.

7 Related Work and Conclusion

There have been a number of OS verification projects, including seL4 [16, 15], Verisoft [4], VCC/VeriSoftXT [9, 3], Verve [29], and CertiKOS [13, 8]. Most of them have no or limited support of preemption and multi-level interrupts.

seL4 [16, 15] is one of the milestone OS kernel verification projects. The verification is fully mechanized in Isabelle/HOL. The kernel of seL4 does not support general preemption. Instead, tasks are preemptible only at specific points. Therefore the code verified is mostly sequential. On the other hand, the seL4 project has verified rich features and properties such as virtual memory, real-time properties and security properties, which are not done in our work.

The Verisoft project also verifies OS microkernels [4] in Isabelle/HOL, but the CVM model used there does not permit interrupts inside the kernel. Its successor project, Verisoft XT [3], uses VCC [9] to verify the commercial Hyper-V

hypervisor. VCC supports verification of concurrent C code by inserting auxiliary code and ghost states. The proofs have a refinement flavor, but VCC does not establish contextual refinement as what we do. Also it is unclear how VCC is applied to verify multi-level nested interrupts in hypervisors.

Verve [29] combines a type-safe kernel with a minimal hardware abstraction layer. The kernel is concurrent, but the properties verified are mostly about type safety, much weaker than our contextual refinement property. Also Verve simply squashes multiple interrupt levels into a single level and does not really handle multi-level interrupts. VCC/VerisoftXT and Verve use the Z3 SMT solver [10] for better automation, while we use Coq which generates machine-checkable proofs. Also the soundness of our program logic is proved in Coq. Therefore the trusted computing base (TCB) of our approach is smaller.

Gu *et al.* [13] verify the mCertiKOS hypervisor. Their kernel is sequential. Recently, Chen *et al.* [8] propose a framework for building certified interruptible OS kernels (based on mCertiKOS) with device drivers. Their framework does not support preemptive concurrency as ours, and it requires that interrupt handlers for device drivers and non-handler kernel code should not share any state.

Gotsman and Yang [12] developed a program logic based on CSL, which decomposes the verification of preemptive kernels into verifying the scheduler and the tasks. Their proofs are on-paper only and not mechanized. The machine model does not support multi-level interrupts, also their program logic is used to prove partial correctness, not contextual refinement as we do.

Conclusion. We have developed a practical verification framework for general verification purpose of preemptive OS kernels with multi-level interrupts. Correctness of the OS kernel is formalized as a contextual refinement between the low-level concrete implementations and the high-level specifications. As far as we know, our work is the first to establish contextual refinement for system APIs of a preemptive OS kernel. We have applied the framework to verify key modules and PIF of $\mu\text{C}/\text{OS-II}$, a commercial embedded real-time OS.

It is worth noting that although our verification framework is developed to verify $\mu\text{C}/\text{OS-II}$, it is a general verification framework and most of its building blocks can be reused to verify other OS kernels. As shown in Fig. 2, the small-step semantics for the C subset, the program logic and the tactics are all general and mostly independent of the $\mu\text{C}/\text{OS-II}$ verification project. A potential limitation is that the interrupt mechanism in our operational semantics is modeled specifically based on the Intel 8259A interrupt controller, and the program logic rules for interrupts are designed accordingly. However, the logic rules follow the general ownership transfer idea from CSL. With a different processor and interrupt mechanism, even though we may need to change the current inference rules for interrupt primitives, we can apply the same ownership transfer idea, and the required change should be superficial. Another limitation is that our C subset is chosen based on the $\mu\text{C}/\text{OS-II}$ code. In particular, it does not allow function pointers, which requires the support of higher-order functions in the logic.

References

- [1] The coq development team: The Coq proof assistant. <http://coq.inria.fr>.
- [2] The real-time kernel: $\mu C/OS-II$. <http://micrium.com/rtos/ucosii/overview>.
- [3] The Verisoft XT Project, 2007. <http://www.verisoftxt.de>.
- [4] E. Alkassar, W. Paul, A. Starostin, and A. Tsyban. Pervasive verification of an OS microkernel: Inline assembly, memory consumption, concurrent devices. In *VSTTE*, pages 71–85, 2010.
- [5] A. W. Appel. Tactics for separation logic, 2006. <http://www.cs.princeton.edu/~appel/papers/septacs.pdf>.
- [6] O. Babaoglu, K. Marzullo, and F. B. Schneider. A formalization of priority inversion. *Real-time Systems*, 5:285–303, 1993.
- [7] J. Cao, M. Fu, and X. Feng. Practical tactics for verifying C programs in coq. In *CPP*, pages 97–108, 2015.
- [8] H. Chen, N. Wu, Z. Shao, J. Lockerman, and R. Gu. Toward compositional verification of interruptible os kernels and device drivers. In *PLDI*, page (to appear), 2016.
- [9] E. Cohen, M. Dahlweid, M. Hillebrand, D. Leinenbach, M. Moskal, T. Santen, W. Schulte, and S. Tobies. VCC: A practical system for verifying concurrent C. In *TPHOLs*, pages 23–42, 2009.
- [10] L. M. de Moura and N. Bjørner. Z3: An efficient SMT solver. In *TACAS*, pages 337–340, 2008.
- [11] X. Feng, Z. Shao, Y. Dong, and Y. Guo. Certifying low-level programs with hardware interrupts and preemptive threads. In *PLDI*, pages 170–182, 2008.
- [12] A. Gotsman and H. Yang. Modular verification of preemptive OS kernels. *J. Funct. Program.*, 23(4):452–514, 2013.
- [13] R. Gu, J. Koenig, T. Ramanananandro, Z. Shao, X. N. Wu, S.-C. Weng, H. Zhang, and Y. Guo. Deep specifications and certified abstraction layers. In *POPL*, pages 595–608, 2015.
- [14] G. Klein, J. Andronick, K. Elphinstone, G. Heiser, D. Cock, P. Derrin, D. Elka-duwe, K. Engelhardt, R. Kolanski, M. Norrish, T. Sewell, H. Tuch, and S. Win-wood. sel4: Formal verification of an operating-system kernel. *Commun. ACM*, 53(6):107–115, 2010.
- [15] G. Klein, J. Andronick, K. Elphinstone, T. C. Murray, T. Sewell, R. Kolanski, and G. Heiser. Comprehensive formal verification of an OS microkernel. *ACM Trans. Comput. Syst.*, 32(1):2, 2014.
- [16] G. Klein, K. Elphinstone, G. Heiser, J. Andronick, D. Cock, P. Derrin, D. Elka-duwe, K. Engelhardt, R. Kolanski, M. Norrish, T. Sewell, H. Tuch, and S. Win-wood. sel4: Formal verification of an os kernel. In *SOSP*, pages 207–220, 2009.
- [17] X. Leroy. Formal verification of a realistic compiler. *Commun. ACM*, 52(7):107–115, 2009.
- [18] H. Liang and X. Feng. Modular verification of linearizability with non-fixed lin-earization points. In *PLDI*, pages 459–470, 2013.
- [19] H. Liang, X. Feng, and M. Fu. A rely-guarantee-based simulation for verifying concurrent program transformations. In *POPL*, pages 455–468, 2012.
- [20] H. Liang, X. Feng, and Z. Shao. Compositional verification of termination-preserving refinement of concurrent programs. In *CSL-LICS*, pages 65:1–65:10, 2014.
- [21] A. McCreight. Practical tactics for separation logic. In *TPHOLs*, pages 343–358, 2009.

- [22] P. W. O’Hearn. Resources, concurrency and local reasoning. In *CONCUR*, pages 49–67, 2004.
- [23] X. Qiu, P. Garg, A. Stefanescu, and P. Madhusudan. Natural proofs for structure, data, and separation. In *ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI ’13, Seattle, WA, USA, June 16-19, 2013*, pages 231–242, 2013.
- [24] J. Sevcík, V. Vafeiadis, F. Z. Nardelli, S. Jagannathan, and P. Sewell. Compccerttso: A verified compiler for relaxed-memory concurrency. *J. ACM*, 60(3):22, 2013.
- [25] L. Sha, R. Rajkumar, and J. P. Lehoczky. Priority inheritance protocols: An approach to real-time synchronization. *IEEE Transactions on Computers*, 39:1175–1185, 1990.
- [26] A. Turon, D. Dreyer, and L. Birkedal. Unifying refinement and hoare-style reasoning in a logic for higher-order concurrency. In *ICFP*, pages 377–390, 2013.
- [27] A. Turon, J. Thamsborg, A. Ahmed, L. Birkedal, and D. Dreyer. Logical relations for fine-grained concurrency. In *POPL*, pages 343–356, 2013.
- [28] F. Xu, M. Fu, X. Feng, X. Zhang, H. Zhang, and Z. Li. A practical verification framework for preemptive OS kernels (technical report and coq implementations), May 2016. <http://staff.ustc.edu.cn/~fuming/research/certiucos>.
- [29] J. Yang and C. Hawblitzel. Safe to the last instruction: automated verification of a type-safe operating system. In *PLDI*, pages 99–110, 2010.