

Priority Inheritance Protocol Proved Correct

Xingyuan Zhang¹, Christian Urban², and Chunhan Wu¹

¹ PLA University of Science and Technology, China

² King's College London, United Kingdom

Abstract. In real-time systems with threads, resource locking and priority scheduling, one faces the problem of Priority Inversion. This problem can make the behaviour of threads unpredictable and the resulting bugs can be hard to find. The Priority Inheritance Protocol is one solution implemented in many systems for solving this problem, but the correctness of this solution has never been formally verified in a theorem prover. As already pointed out in the literature, the original informal investigation of the Property Inheritance Protocol presents a correctness “proof” for an *incorrect* algorithm. In this paper we fix the problem of this proof by making all notions precise and implementing a variant of a solution proposed earlier. Our formalisation in Isabelle/HOL uncovers facts not mentioned in the literature, but also shows how to efficiently implement this protocol. Earlier correct implementations were criticised as too inefficient. Our formalisation is based on Paulson’s inductive approach to verifying protocols.

Keywords: Priority Inheritance Protocol, formal correctness proof, real-time systems, Isabelle/HOL

1 Introduction

Many real-time systems need to support threads involving priorities and locking of resources. Locking of resources ensures mutual exclusion when accessing shared data or devices that cannot be preempted. Priorities allow scheduling of threads that need to finish their work within deadlines. Unfortunately, both features can interact in subtle ways leading to a problem, called *Priority Inversion*. Suppose three threads having priorities H (igh), M (edium) and L (ow). We would expect that the thread H blocks any other thread with lower priority and itself cannot be blocked by any thread with lower priority. Alas, in a naive implementation of resource locking and priorities this property can be violated. Even worse, H can be delayed indefinitely by threads with lower priorities. For this let L be in the possession of a lock for a resource that also H needs. H must therefore wait for L to exit the critical section and release this lock. The problem is that L might in turn be blocked by any thread with priority M , and so H sits there potentially waiting indefinitely. Since H is blocked by threads with lower priorities, the problem is called Priority Inversion. It was first described in [4] in the context of the Mesa programming language designed for concurrent programming.

If the problem of Priority Inversion is ignored, real-time systems can become unpredictable and resulting bugs can be hard to diagnose. The classic example where this happened is the software that controlled the Mars Pathfinder mission in 1997 [6]. Once the spacecraft landed, the software shut down at irregular intervals leading to loss of project time as normal operation of the craft could only resume the next day (the mission and data already collected were fortunately not lost, because of a clever system design). The reason for the shutdowns was that the scheduling software fell victim of Priority Inversion: a low priority thread locking a resource prevented a high priority thread from running in time leading to a system reset. Once the problem was found, it was rectified by enabling the *Priority Inheritance Protocol* (PIP) [7]³ in the scheduling software.

The idea behind PIP is to let the thread L temporarily inherit the high priority from H until L leaves the critical section unlocking the resource. This solves the problem of H having to wait indefinitely, because L cannot be blocked by threads having priority M . While a few other solutions exist for the Priority Inversion problem, PIP is one that is widely deployed and implemented. This includes VxWorks (a proprietary real-time OS used in the Mars Pathfinder mission, in Boeing’s 787 Dreamliner, Honda’s ASIMO robot, etc.), but also the POSIX 1003.1c Standard realised for example in libraries for FreeBSD, Solaris and Linux.

One advantage of PIP is that increasing the priority of a thread can be dynamically calculated by the scheduler. This is in contrast to, for example, *Priority Ceiling* [7], another solution to the Priority Inversion problem, which requires static analysis of the program in order to prevent Priority Inversion. However, there has also been strong criticism against PIP. For instance, PIP cannot prevent deadlocks when lock dependencies are circular, and also blocking times can be substantial (more than just the duration of a critical section). Though, most criticism against PIP centres around unreliable implementations and PIP being too complicated and too inefficient. For example, Yodaiken writes in [12]:

“Priority inheritance is neither efficient nor reliable. Implementations are either incomplete (and unreliable) or surprisingly complex and intrusive.”

He suggests to avoid PIP altogether by not allowing critical sections to be preempted. Unfortunately, this solution does not help in real-time systems with hard deadlines for high-priority threads.

In our opinion, there is clearly a need for investigating correct algorithms for PIP. A few specifications for PIP exist (in English) and also a few high-level descriptions of implementations (e.g. in the textbook [9, Section 5.6.5]), but they help little with actual implementations. That this is a problem in practise is proved by an email from Baker, who wrote on 13 July 2009 on the Linux Kernel mailing list:

“I observed in the kernel code (to my disgust), the Linux PIP implementation is a nightmare: extremely heavy weight, involving maintenance of a full wait-for graph, and requiring updates for a range of events, including priority changes and interruptions of wait operations.”

³ Sha et al. call it the *Basic Priority Inheritance Protocol* [7] and others sometimes also call it *Priority Boosting*.

The criticism by Yodaiken, Baker and others suggests to us to look again at PIP from a more abstract level (but still concrete enough to inform an implementation), and makes PIP an ideal candidate for a formal verification. One reason, of course, is that the original presentation of PIP [7], despite being informally “proved” correct, is actually *flawed*.

Yodaiken [12] points to a subtlety that had been overlooked in the informal proof by Sha et al. They specify in [7] that after the thread (whose priority has been raised) completes its critical section and releases the lock, it “returns to its original priority level.” This leads them to believe that an implementation of PIP is “rather straightforward” [7]. Unfortunately, as Yodaiken points out, this behaviour is too simplistic. Consider the case where the low priority thread L locks *two* resources, and two high-priority threads H and H' each wait for one of them. If L releases one resource so that H , say, can proceed, then we still have Priority Inversion with H' (which waits for the other resource). The correct behaviour for L is to revert to the highest remaining priority of the threads that it blocks. The advantage of formalising the correctness of a high-level specification of PIP in a theorem prover is that such issues clearly show up and cannot be overlooked as in informal reasoning (since we have to analyse all possible behaviours of threads, i.e. *traces*, that could possibly happen).

Contributions: There have been earlier formal investigations into PIP [2,3,11], but they employ model checking techniques. This paper presents a formalised and mechanically checked proof for the correctness of PIP (to our knowledge the first one; the earlier informal proof by Sha et al. [7] is flawed). In contrast to model checking, our formalisation provides insight into why PIP is correct and allows us to prove stronger properties that, as we will show, can inform an efficient implementation. For example, we found by “playing” with the formalisation that the choice of the next thread to take over a lock when a resource is released is irrelevant for PIP being correct. Something which has not been mentioned in the relevant literature.

2 Formal Model of the Priority Inheritance Protocol

The Priority Inheritance Protocol, short PIP, is a scheduling algorithm for a single-processor system.⁴ Our model of PIP is based on Paulson’s inductive approach to protocol verification [5], where the *state* of a system is given by a list of events that happened so far. *Events* of PIP fall into five categories defined as the datatype:

datatype <i>event</i>	=	<i>Create thread priority</i>	
		<i>Exit thread</i>	
		<i>Set thread priority</i>	reset of the priority for <i>thread</i>
		<i>P thread cs</i>	request of resource <i>cs</i> by <i>thread</i>
		<i>V thread cs</i>	release of resource <i>cs</i> by <i>thread</i>

whereby threads, priorities and (critical) resources are represented as natural numbers. The event *Set* models the situation that a thread obtains a new priority given by the

⁴ We shall come back later to the case of PIP on multi-processor systems.

programmer or user (for example via the `nice` utility under UNIX). As in Paulson’s work, we need to define functions that allow us to make some observations about states. One, called *threads*, calculates the set of “live” threads that we have seen so far:

$$\begin{aligned} \text{threads } [] & \stackrel{\text{def}}{=} \emptyset \\ \text{threads } (\text{Create } th \text{ prio}::s) & \stackrel{\text{def}}{=} \{th\} \cup \text{threads } s \\ \text{threads } (\text{Exit } th::s) & \stackrel{\text{def}}{=} \text{threads } s - \{th\} \\ \text{threads } (::_s) & \stackrel{\text{def}}{=} \text{threads } s \end{aligned}$$

In this definition `::_` stands for list-cons. Another function calculates the priority for a thread *th*, which is defined as

$$\begin{aligned} \text{priority } th [] & \stackrel{\text{def}}{=} 0 \\ \text{priority } th (\text{Create } th' \text{ prio}::s) & \stackrel{\text{def}}{=} \text{if } th' = th \text{ then } prio \text{ else } \text{priority } th s \\ \text{priority } th (\text{Set } th' \text{ prio}::s) & \stackrel{\text{def}}{=} \text{if } th' = th \text{ then } prio \text{ else } \text{priority } th s \\ \text{priority } th (::_s) & \stackrel{\text{def}}{=} \text{priority } th s \end{aligned}$$

In this definition we set *0* as the default priority for threads that have not (yet) been created. The last function we need calculates the “time”, or index, at which time a process had its priority last set.

$$\begin{aligned} \text{last_set } th [] & \stackrel{\text{def}}{=} 0 \\ \text{last_set } th (\text{Create } th' \text{ prio}::s) & \stackrel{\text{def}}{=} \text{if } th = th' \text{ then } |s| \text{ else } \text{last_set } th s \\ \text{last_set } th (\text{Set } th' \text{ prio}::s) & \stackrel{\text{def}}{=} \text{if } th = th' \text{ then } |s| \text{ else } \text{last_set } th s \\ \text{last_set } th (::_s) & \stackrel{\text{def}}{=} \text{last_set } th s \end{aligned}$$

In this definition `|s|` stands for the length of the list of events *s*. Again the default value in this function is *0* for threads that have not been created yet. A *precedence* of a thread *th* in a state *s* is the pair of natural numbers defined as

$$\text{prec } th s \stackrel{\text{def}}{=} (\text{priority } th s, \text{last_set } th s)$$

The point of precedences is to schedule threads not according to priorities (because what should we do in case two threads have the same priority), but according to precedences. Precedences allow us to always discriminate between two threads with equal priority by taking into account the time when the priority was last set. We order precedences so that threads with the same priority get a higher precedence if their priority has been set earlier, since for such threads it is more urgent to finish their work. In an implementation this choice would translate to a quite natural FIFO-scheduling of processes with the same priority.

Next, we introduce the concept of *waiting queues*. They are lists of threads associated with every resource. The first thread in this list (i.e. the head, or short *hd*) is chosen to be the one that is in possession of the “lock” of the corresponding resource. We model waiting queues as functions, below abbreviated as *wq*. They take a resource as argument and return a list of threads. This allows us to define when a thread *holds*, respectively *waits* for, a resource *cs* given a waiting queue function *wq*.

$$\begin{aligned} \text{holds } wq \ th \ cs &\stackrel{\text{def}}{=} th \in \text{set} (wq \ cs) \wedge th = hd (wq \ cs) \\ \text{waits } wq \ th \ cs &\stackrel{\text{def}}{=} th \in \text{set} (wq \ cs) \wedge th \neq hd (wq \ cs) \end{aligned}$$

In this definition we assume *set* converts a list into a set. At the beginning, that is in the state where no thread is created yet, the waiting queue function will be the function that returns the empty list for every resource.

$$\text{all_unlocked} \stackrel{\text{def}}{=} \lambda _. [] \quad (1)$$

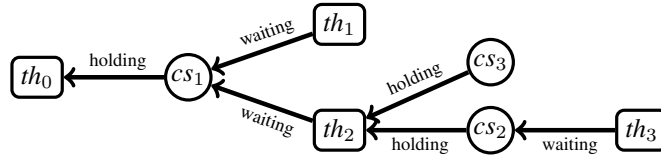
Using *holds* and *waits*, we can introduce *Resource Allocation Graphs* (RAG), which represent the dependencies between threads and resources. We represent RAGs as relations using pairs of the form

$$(T \ th, \ C \ cs) \quad \text{and} \quad (C \ cs, \ T \ th)$$

where the first stands for a *waiting edge* and the second for a *holding edge* (*C* and *T* are constructors of a datatype for vertices). Given a waiting queue function, a RAG is defined as the union of the sets of waiting and holding edges, namely

$$\text{RAG } wq \stackrel{\text{def}}{=} \{(T \ th, \ C \ cs) \mid \text{waits } wq \ th \ cs\} \cup \{(C \ cs, \ T \ th) \mid \text{holds } wq \ th \ cs\}$$

Given three threads and three resources, an instance of a RAG can be pictured as follows:



The use of relations for representing RAGs allows us to conveniently define the notion of the *dependants* of a thread using the transitive closure operation for relations. This gives

$$\text{dependants } wq \ th \stackrel{\text{def}}{=} \{th' \mid (T \ th', \ T \ th) \in (\text{RAG } wq)^+\}$$

This definition needs to account for all threads that wait for a thread to release a resource. This means we need to include threads that transitively wait for a resource being released (in the picture above this means the dependants of *th*₀ are *th*₁ and *th*₂, which wait for resource *cs*₁, but also *th*₃, which cannot make any progress unless *th*₂ makes progress, which in turn needs to wait for *th*₀ to finish). If there is a circle in a RAG, then clearly we have a deadlock. Therefore when a thread requests a resource, we must ensure that the resulting RAG is not circular.

Next we introduce the notion of the *current precedence* of a thread *th* in a state *s*. It is defined as

$$\text{cprec } wq \ s \ th \stackrel{\text{def}}{=} \text{Max} (\{\text{prec } th \ s\} \cup \{\text{prec } th' \ s \mid th' \in \text{dependants } wq \ th\}) \quad (2)$$

where the dependants of th are given by the waiting queue function. While the precedence $prec$ of a thread is determined by the programmer (for example when the thread is created), the point of the current precedence is to let the scheduler increase this precedence, if needed according to PIP. Therefore the current precedence of th is given as the maximum of the precedence th has in state s and all threads that are dependants of th . Since the notion *dependants* is defined as the transitive closure of all dependent threads, we deal correctly with the problem in the informal algorithm by Sha et al. [7] where a priority of a thread is lowered prematurely.

The next function, called *schs*, defines the behaviour of the scheduler. It will be defined by recursion on the state (a list of events); this function returns a *schedule state*, which we represent as a record consisting of two functions:

$$(\text{wq_fun}, \text{cprec_fun})$$

The first function is a waiting queue function (that is, it takes a resource cs and returns the corresponding list of threads that lock, respectively wait for, it); the second is a function that takes a thread and returns its current precedence (see (2)). We assume the usual getter and setter methods for such records.

In the initial state, the scheduler starts with all resources unlocked (the corresponding function is defined in (1)) and the current precedence of every thread is initialised with $(0, 0)$; that means $\text{initial_cprec} \stackrel{\text{def}}{=} \lambda _ . (0, 0)$. Therefore we have for the initial state

$$\text{schs} [] \stackrel{\text{def}}{=} (\text{wq_fun} = \text{all_unlocked}, \text{cprec_fun} = \text{initial_cprec})$$

The cases for *Create*, *Exit* and *Set* are also straightforward: we calculate the waiting queue function of the (previous) state s ; this waiting queue function wq is unchanged in the next schedule state—because none of these events lock or release any resource; for calculating the next cprec_fun , we use wq and cprec . This gives the following three clauses for *schs*:

$$\begin{aligned} \text{schs} (\text{Create } th \text{ prio}::s) &\stackrel{\text{def}}{=} \\ &\text{let } wq = \text{wq_fun} (\text{schs } s) \text{ in} \\ &(\text{wq_fun} = wq, \text{cprec_fun} = \text{cprec } wq (\text{Create } th \text{ prio}::s)) \\ \text{schs} (\text{Exit } th::s) &\stackrel{\text{def}}{=} \\ &\text{let } wq = \text{wq_fun} (\text{schs } s) \text{ in} \\ &(\text{wq_fun} = wq, \text{cprec_fun} = \text{cprec } wq (\text{Exit } th::s)) \\ \text{schs} (\text{Set } th \text{ prio}::s) &\stackrel{\text{def}}{=} \\ &\text{let } wq = \text{wq_fun} (\text{schs } s) \text{ in} \\ &(\text{wq_fun} = wq, \text{cprec_fun} = \text{cprec } wq (\text{Set } th \text{ prio}::s)) \end{aligned}$$

More interesting are the cases where a resource, say cs , is locked or released. In these cases we need to calculate a new waiting queue function. For the event $P \text{ th } cs$, we have to update the function so that the new thread list for cs is the old thread list plus the thread th appended to the end of that list (remember the head of this list is assigned to be in the possession of this resource). This gives the clause

$$\begin{aligned}
 schs (P th cs::s) &\stackrel{def}{=} \\
 &\text{let } wq = wq_fun (schs s) \text{ in} \\
 &\text{let } new_wq = wq(cs := (wq cs @ [th])) \text{ in} \\
 &(\!|wq_fun = new_wq, cprec_fun = cprec new_wq (P th cs::s)\!|)
 \end{aligned}$$

The clause for event $V th cs$ is similar, except that we need to update the waiting queue function so that the thread that possessed the lock is deleted from the corresponding thread list. For this list transformation, we use the auxiliary function *release*. A simple version of *release* would just delete this thread and return the remaining threads, namely

$$\begin{aligned}
 release [] &\stackrel{def}{=} [] \\
 release (_::qs) &\stackrel{def}{=} qs
 \end{aligned}$$

In practice, however, often the thread with the highest precedence in the list will get the lock next. We have implemented this choice, but later found out that the choice of which thread is chosen next is actually irrelevant for the correctness of PIP. Therefore we prove the stronger result where *release* is defined as

$$\begin{aligned}
 release [] &\stackrel{def}{=} [] \\
 release (_::qs) &\stackrel{def}{=} SOME qs'. distinct qs' \wedge set qs' = set qs
 \end{aligned}$$

where *SOME* stands for Hilbert's epsilon and implements an arbitrary choice for the next waiting list. It just has to be a list of distinctive threads and contain the same elements as *qs*. This gives for V the clause:

$$\begin{aligned}
 schs (V th cs::s) &\stackrel{def}{=} \\
 &\text{let } wq = wq_fun (schs s) \text{ in} \\
 &\text{let } new_wq = release (wq cs) \text{ in} \\
 &(\!|wq_fun = new_wq, cprec_fun = cprec new_wq (V th cs::s)\!|)
 \end{aligned}$$

Having the scheduler function *schs* at our disposal, we can “lift”, or overload, the notions *waits*, *holds*, *RAG* and *cprec* to operate on states only.

$$\begin{aligned}
 holds s &\stackrel{def}{=} holds (wq_fun (schs s)) \\
 waits s &\stackrel{def}{=} waits (wq_fun (schs s)) \\
 RAG s &\stackrel{def}{=} RAG (wq_fun (schs s)) \\
 cprec s &\stackrel{def}{=} cprec_fun (schs s)
 \end{aligned}$$

With these abbreviations we can introduce the notion of threads being *ready* in a state (i.e. threads that do not wait for any resource) and the running thread.

$$\begin{aligned}
 ready s &\stackrel{def}{=} \{th \in threads s \mid \forall cs. \neg waits s th cs\} \\
 running s &\stackrel{def}{=} \{th \in ready s \mid cprec s th = Max (cprec s ' ready s)\}
 \end{aligned}$$

In this definition $_ _$ stands for the image of a set under a function. Note that in the initial state, that is where the list of events is empty, the set *threads* is empty and therefore there is neither a thread ready nor running. If there is one or more threads ready, then there can only be *one* thread running, namely the one whose current precedence is equal to the maximum of all ready threads. We use sets to capture both possibilities. We can now also conveniently define the set of resources that are locked by a thread in a given state.

$$resources\ s\ th \stackrel{def}{=} \{cs \mid (C\ cs, T\ th) \in RAG\ s\}$$

Finally we can define what a *valid state* is in our model of PIP. For example we cannot expect to be able to exit a thread, if it was not created yet. These validity constraints on states are characterised by the inductive predicate *step* and *valid_state*. We first give five inference rules for *step* relating a state and an event that can happen next.

$$\frac{th \notin threads\ s}{step\ s\ (Create\ th\ prio)} \qquad \frac{th \in running\ s \quad resources\ s\ th = \emptyset}{step\ s\ (Exit\ th)}$$

The first rule states that a thread can only be created, if it does not yet exist. Similarly, the second rule states that a thread can only be terminated if it was running and does not lock any resources anymore (this simplifies slightly our model; in practice we would expect the operating system releases all locks held by a thread that is about to exit). The event *Set* can happen if the corresponding thread is running.

$$\frac{th \in running\ s}{step\ s\ (Set\ th\ prio)}$$

If a thread wants to lock a resource, then the thread needs to be running and also we have to make sure that the resource lock does not lead to a cycle in the RAG. In practice, ensuring the latter is the responsibility of the programmer. In our formal model we brush aside these problematic cases in order to be able to make some meaningful statements about PIP.⁵

$$\frac{th \in running\ s \quad (C\ cs, T\ th) \notin (RAG\ s)^+}{step\ s\ (P\ th\ cs)}$$

Similarly, if a thread wants to release a lock on a resource, then it must be running and in the possession of that lock. This is formally given by the last inference rule of *step*.

$$\frac{th \in running\ s \quad holds\ s\ th\ cs}{step\ s\ (V\ th\ cs)}$$

A valid state of PIP can then be conveniently be defined as follows:

⁵ This situation is similar to the infamous occurs check in Prolog: In order to say anything meaningful about unification, one needs to perform an occurs check. But in practice the occurs check is omitted and the responsibility for avoiding problems rests with the programmer.

$$\frac{}{\text{valid_state } \square} \quad \frac{\text{valid_state } s \quad \text{step } s \ e}{\text{valid_state } (e::s)}$$

This completes our formal model of PIP. In the next section we present properties that show our model of PIP is correct.

3 The Correctness Proof

Sha et al. [7, Theorem 6] state their correctness criterion for PIP in terms of the number of critical resources: if there are m critical resources, then a blocked job with high priority can only be blocked m times—that is a *bounded* number of times. This result on its own, strictly speaking, does *not* prevent indefinite, or unbounded, Priority Inversion, because if one low-priority thread does not give up its critical resource (the one the high-priority thread is waiting for), then the high-priority thread can never run. The argument of Sha et al. is that *if* threads release locked resources in a finite amount of time, then indefinite Priority Inversion cannot occur—the high-priority thread is guaranteed to run eventually. The assumption is that programmers always ensure that this is the case. However, even taking this assumption into account, their correctness property is *not* true for their version of PIP. As Yodaiken [12] pointed out: If a low-priority thread possesses locks to two resources for which two high-priority threads are waiting for, then lowering the priority prematurely after giving up only one lock, can cause indefinite Priority Inversion for one of the high-priority threads, invalidating their bound.

Even when fixed, their proof idea does not seem to go through for us, because of the way we have set up our formal model of PIP. The reason is that we allow that critical sections can intersect (something Sha et al. explicitly exclude). Therefore we have a different correctness criterion for PIP. The idea behind our criterion is as follows: for all states s , we know the corresponding thread th with the highest precedence; we show that in every future state (denoted by $s' @ s$) in which th is still alive, either th is running or it is blocked by a thread that was alive in the state s . Since in s , as in every state, the set of alive threads is finite, th can only be blocked a finite number of times. We will actually prove a stricter bound below. However, this correctness criterion hinges upon a number of assumptions about the states s and $s' @ s$, the thread th and the events happening in s' . We list them next:

Assumptions on the states s and $s' @ s$: In order to make any meaningful statement, we need to require that s and $s' @ s$ are valid states, namely

$$\text{valid_state } s \\ \text{valid_state } (s' @ s)$$

Assumptions on the thread th : The thread th must be alive in s and has the highest precedence of all alive threads in s . Furthermore the priority of th is $prio$ (we need this in the next assumptions).

$$th \in \text{threads } s \\ \text{prec } th \ s = \text{Max } (cprec \ s \ \text{threads } s) \\ \text{prec } th \ s = (prio, -)$$

Assumptions on the events in s' : We want to prove that th cannot be blocked indefinitely. Of course this can happen if threads with higher priority than th are continuously created in s' . Therefore we have to assume that events in s' can only create (respectively set) threads with equal or lower priority than $prio$ of th . We also need to assume that the priority of th does not get reset and also that th does not get “exited” in s' . This can be ensured by assuming the following three implications.

If Create $th' prio' \in set\ s'$ then $prio' \leq prio$
If Set $th' prio' \in set\ s'$ then $th' \neq th$ and $prio' \leq prio$
If Exit $th' \in set\ s'$ then $th' \neq th$

Under these assumptions we will prove the following correctness property:

Theorem 1. *Given the assumptions about states s and $s' @ s$, the thread th and the events in s' , if $th' \in running\ (s' @ s)$ and $th' \neq th$ then $th' \in threads\ s$.*

This theorem ensures that the thread th , which has the highest precedence in the state s , can only be blocked in the state $s' @ s$ by a thread th' that already existed in s . As we shall see shortly, that means by only finitely many threads. Like in the argument by Sha et al. this finite bound does not guarantee absence of indefinite Priority Inversion. For this we further have to assume that every thread gives up its resources after a finite amount of time. We found that this assumption is awkward to formalise in our model. Therefore we leave it out and let the programmer assume the responsibility to program threads in such a benign manner (in addition to causing no circularity in the RAG). In this detail, we do not make any progress in comparison with the work by Sha et al.

In what follows we will describe properties of PIP that allow us to prove Theorem 1 and, when instructive, briefly describe our argument. It is relatively easily to see that

running $s \subseteq ready\ s \subseteq threads\ s$
If valid_state s then finite (threads s).

whereby the second property is by induction of *valid_state*. The next three properties are

If valid_state s and waits $s\ th\ cs_1$ and waits $s\ th\ cs_2$ then $cs_1 = cs_2$.
If holds $s\ th_1\ cs$ and holds $s\ th_2\ cs$ then $th_1 = th_2$.
If valid_state s and $th_1 \in running\ s$ and $th_2 \in running\ s$ then $th_1 = th_2$.

The first property states that every waiting thread can only wait for a single resource (because it gets suspended after requesting that resource); the second that every resource can only be held by a single thread; the third property establishes that in every given valid state, there is at most one running thread. We can also show the following properties about the RAG in s .

If valid_state s then:
acyclic (RAG s), finite (RAG s) and wf ((RAG s)⁻¹),
if $T\ th \in Domain\ (RAG\ s)$ then $th \in threads\ s$ and
if $T\ th \in Range\ (RAG\ s)$ then $th \in threads\ s$.

The acyclicity property follow from how we restricted the events in *step*; similarly the finiteness and well-foundedness property. The last two properties establish that every thread in a RAG (either holding or waiting for a resource) is a live thread.

To state the key lemma in our proof, it will be convenient to introduce the notion of a *detached* thread in a state, that is one which does not hold any critical resource nor requests one.

Lemma 1. *Given the assumptions about states s and $s' @ s$, the thread th and the events in s' , if $th' \in \text{treads}(s' @ s)$, $th' \neq th$ and $\text{detached}(s' @ s) th'$ then $th' \notin \text{running}(s' @ s)$.*

The point of this lemma is that a thread different from th (which has the highest precedence in s) not holding any resource cannot be running in the state $s' @ s$.

Proof. Since thread th' does not hold any resource, no thread can depend on it. Therefore its current precedence $cprec(s' @ s) th'$ equals its own precedence $prec th'(s' @ s)$. Since th has the highest precedence in the state $(s' @ s)$ and precedences are distinct among threads, we have $prec th'(s' @ s) < prec th(s' @ s)$. From this we have $cprec(s' @ s) th' < prec th(s' @ s)$. Since $prec th(s' @ s)$ is already the highest $cprec(s' @ s) th$ can not be higher than this and can not be lower either (by definition of $cprec$). Consequently, we have $prec th(s' @ s) = cprec(s' @ s) th$. Finally we have $cprec(s' @ s) th' < cprec(s' @ s) th$. By definition of *running*, th' can not be running in state $s' @ s$, as we had to show. \square

Since th' is not able to run at state $s' @ s$, it is not able to issue a text "P" or V event. Therefore if $s' @ s$ is extended one step further, th' still cannot hold any resource. The situation will not change in further extensions as long as th holds the highest precedence.

The following lemmas show how every node in RAG can be chased to ready threads:

1. Every node in RAG can be chased to a ready thread (*chain_building*):

$$\begin{aligned} & \llbracket \text{valid_state } s; \text{ node} \in \text{Domain}(RAG\ s) \rrbracket \\ & \implies \exists th'. th' \in \text{ready } s \wedge (\text{node}, T\ th') \in (RAG\ s)^+ \end{aligned}$$

2. The ready thread chased to is unique (*dchain_unique*):

$$\begin{aligned} & \llbracket \text{valid_state } s; (n, T\ th_1) \in (RAG\ s)^+; th_1 \in \text{ready } s; (n, T\ th_2) \in (RAG\ s)^+; \\ & th_2 \in \text{ready } s \rrbracket \\ & \implies th_1 = th_2 \end{aligned}$$

Some deeper results about the system:

1. The maximum of $cprec$ and $prec$ are equal (*max_cp_eq*):

$$\begin{aligned} & \text{valid_state } s \implies \\ & \text{Max}(cprec\ s \text{ ' threads } s) = \text{Max}((\lambda th. prec\ th\ s) \text{ ' threads } s) \end{aligned}$$

2. There must be one ready thread having the max $cprec$ -value (*max_cp_readys_threads*):

$$\text{valid_state } s \implies \text{Max } (cprec \ s \ ' \ ready \ s) = \text{Max } (cprec \ s \ ' \ threads \ s)$$

The relationship between the count of P and V and the number of critical resources held by a thread is given as follows:

1. The V -operation decreases the number of critical resources one thread holds (*cntCS_v_dec*)

$$\begin{aligned} \text{valid_state } (V \ \text{thread } cs::s) &\implies \\ \text{cntCS } (V \ \text{thread } cs::s) \ \text{thread} + 1 &= \text{cntCS } s \ \text{thread} \end{aligned}$$

2. The number of V never exceeds the number of P (*cnp_cnv_cncls*):

$$\begin{aligned} \text{valid_state } s &\implies \\ \text{cntP } s \ \text{th} &= \\ \text{cntV } s \ \text{th} + & \\ (\text{if } \text{th} \in \text{ready } s \vee \text{th} \notin \text{threads } s &\text{ then } \text{cntCS } s \ \text{th} \text{ else } \text{cntCS } s \ \text{th} + 1) \end{aligned}$$

3. The number of V equals the number of P when the relevant thread is not living: (*cnp_cnv_eq*):

$$\llbracket \text{valid_state } s; \text{th} \notin \text{threads } s \rrbracket \implies \text{detached } s \ \text{th}$$

4. When a thread is not living, it does not hold any critical resource (*not_thread_holdents*):

$$\llbracket \text{valid_state } s; \text{th} \notin \text{threads } s \rrbracket \implies \text{resources } s \ \text{th} = \emptyset$$

5. When the number of P equals the number of V , the relevant thread does not hold any critical resource, therefore no thread can depend on it (*count_eq_dependents*):

$$\llbracket \text{valid_state } s; \text{detached } s \ \text{th} \rrbracket \implies \text{dependants } (wq \ s) \ \text{th} = \emptyset$$

$$\text{running } (t \ @ \ s) \neq \emptyset$$

3.1 Proof idea

The reason that only threads which already held some resources can be running and block th is that if, otherwise, one thread does not hold any resource, it may never have its priority raised and will not get a chance to run. This fact is supported by lemma *moment_blocked*:

$$\begin{aligned} \llbracket \text{th}' \neq \text{th}; \text{th}' \in \text{threads } (\text{moment } i \ t \ @ \ s); \text{detached } (\text{moment } i \ t \ @ \ s) \ \text{th}'; \\ i \leq j \rrbracket \\ \implies \text{detached } (\text{moment } j \ t \ @ \ s) \ \text{th}' \wedge \\ \text{th}' \in \text{threads } (\text{moment } j \ t \ @ \ s) \wedge \text{th}' \notin \text{running } (\text{moment } j \ t \ @ \ s) \end{aligned}$$

When instantiating i to 0 , the lemma means threads which did not hold any resource in state s will not have a change to run later. Rephrased, it means any thread which is running after th became the highest must have already held some resource at state s .

When instantiating i to a number larger than 0 , the lemma means if a thread releases all its resources at some moment in t , after that, it may never get a change to run. If every thread releases its resource in finite duration, then after a while, only thread th is left running. This shows how indefinite priority inversion can be avoided.

So, the key of the proof is to establish the correctness of *moment_blocked*. We are going to show how this lemma is proved. At the heart of this proof, is lemma *pv_blocked*:

$$\begin{aligned} & \llbracket th' \in \text{threads } (t @ s); th' \neq th; \text{detached } (t @ s) th \rrbracket \\ & \implies th' \notin \text{running } (t @ s) \end{aligned}$$

This lemma says: for any s -extension text " t ", if thread th' does not hold any resource, it can not be running at $t@s$.

Proof:

1. Since thread th' does not hold any resource, no thread may depend on it, so its current precedence $cp (t@s) th'$ equals to its own precedence $preced th' (t@s)$.
2. Since th has the highest precedence in the system and precedences are distinct among threads, we have $preced th' (t@s) < preced th (t@s)$. From this and item 1, we have $cp (t@s) th' < preced th (t@s)$.
3. Since $preced th (t@s)$ is already the highest in the system, $cp (t@s) th$ can not be higher than this and can not be lower neither (by the definition of cp), we have $preced th (t@s) = cp (t@s) th$.
4. Finally we have $cp (t@s) th' < cp (t@s) th$.
5. By definition of *running*, th' can not be running at $t@s$.

Since th' is not able to run at state $t@s$, it is not able to make either text "P" or V action, so if $t@s$ is extended one step further, th' still does not hold any resource. The situation will not unchanged in further extensions as long as th holds the highest precedence. Since this t is arbitrarily chosen except being constrained by predicate *extend_highest_gen* and this predicate has the property that if it holds for t , it also holds for any moment i inside t , as shown by lemma *red_moment*:

$$\begin{aligned} & \text{extend_highest_gen } s \text{ th prio } tm \ t \implies \\ & \text{extend_highest_gen } s \text{ th prio } tm \ (\text{moment } i \ t) \end{aligned}$$

so *pv_blocked* can be applied to any *moment* i t . From this, lemma *moment_blocked* follows.

4 Key properties

All these assumptions are put into a predicate *extend_highest_gen*. It can be proved that *extend_highest_gen* holds for any moment i in it t (*red_moment*):

$$\text{extend_highest_gen } s \text{ th prio } tm \ (\text{moment } i \ t)$$

From this, an induction principle can be derived for t , so that properties already derived for t can be applied to any prefix of t in the proof of new properties about t (*ind*):

$$\begin{aligned} & \llbracket R \rrbracket; \\ & \bigwedge e t. \llbracket \text{valid_state } (t @ s); \text{step } (t @ s) e; \\ & \quad \text{extend_highest_gen } s \text{ th prio } tm t; \\ & \quad \text{extend_highest_gen } s \text{ th prio } tm (e::t); R t \rrbracket \\ & \implies R (e::t) \rrbracket \\ \implies R t \end{aligned}$$

The following properties can be proved about th in t :

1. In t , thread th is kept live and its precedence is preserved as well (th_kept):

$$th \in \text{threads } (t @ s) \wedge \text{prec } th (t @ s) = \text{prec } th s$$
2. In t , thread th 's precedence is always the maximum among all living threads (max_preced):

$$\text{prec } th (t @ s) = \text{Max } ((\lambda th'. \text{prec } th' (t @ s)) \text{ ' } \text{threads } (t @ s))$$
3. In t , thread th 's current precedence is always the maximum precedence among all living threads ($th_cp_max_preced$):

$$cprec (t @ s) th = \text{Max } ((\lambda th'. \text{prec } th' (t @ s)) \text{ ' } \text{threads } (t @ s))$$
4. In t , thread th 's current precedence is always the maximum current precedence among all living threads (th_cp_max):

$$cprec (t @ s) th = \text{Max } (cprec (t @ s) \text{ ' } \text{threads } (t @ s))$$
5. In t , thread th 's current precedence equals its precedence at moment s (th_cp_preced):

$$cprec (t @ s) th = \text{prec } th s$$

The main theorem of this part is to characterizing the running thread during t ($running_inversion_2$):

$$\begin{aligned} th' \in \text{running } (t @ s) & \implies \\ th' = th \vee th' \neq th \wedge th' \in \text{threads } s \wedge cntV s th' < cntP s th' \end{aligned}$$

According to this, if a thread is running, it is either th or was already live and held some resource at moment s (expressed by: $cntV s th' < cntP s th'$).

Since there are only finite many threads live and holding some resource at any moment, if every such thread can release all its resources in finite duration, then after finite duration, none of them may block th anymore. So, no priority inversion may happen then.

5 Properties for an Implementation

While a formal correctness proof for our model of PIP is certainly attractive (especially in light of the flawed proof by Sha et al. [7]), we found that the formalisation can even help us with efficiently implementing PIP.

For example Baker complained that calculating the current precedence in PIP is quite “heavy weight” in Linux (see the Introduction). In our model of PIP the current precedence of a thread in a state s depends on all its dependants—a “global” transitive notion, which is indeed heavy weight (see Def. shown in (2)). We can however improve upon this. For this let us define the notion of *children* of a thread th in a state s as

$$children\ s\ th \stackrel{def}{=} \{th' \mid \exists cs. (T\ th', C\ cs) \in RAG\ s \wedge (C\ cs, T\ th) \in RAG\ s\}$$

where a child is a thread that is one ‘‘hop’’ away from the thread th in the RAG (and waiting for th to release a resource). We can prove that

Lemma 2. *If valid_state s then*

$$cprec\ s\ th = Max(\{prec\ th\ s\} \cup cprec\ s'\ children\ s\ th).$$

That means the current precedence of a thread th can be computed locally by considering only the children of th . In effect, it only needs to be recomputed for th when one of its children changes its current precedence. Once the current precedence is computed in this more efficient manner, the selection of the thread with highest precedence from a set of ready threads is a standard scheduling operation implemented in most operating systems.

Of course the main implementation work for PIP involves the scheduler and coding how it should react to events. Below we outline how our formalisation guides this implementation for each kind of event.

Create th prio: We assume that the current state s' and the next state $s \stackrel{def}{=} Create\ th\ prio::s'$ are both valid (meaning the event is allowed to occur). In this situation we can show that

$$\begin{aligned} RAG\ s &= RAG\ s', \\ cprec\ s\ th &= prec\ th\ s, \text{ and} \\ \text{If } th' \neq th &\text{ then } cprec\ s\ th' = cprec\ s'\ th'. \end{aligned}$$

This means we do not have to recalculate the RAG and also none of the current precedences of the other threads. The current precedence of the created thread th is just its precedence, namely the pair $(prio, |s|)$.

Exit th : We again assume that the current state s' and the next state $s \stackrel{def}{=} Exit\ th::s'$ are both valid. We can show that

$$\begin{aligned} RAG\ s &= RAG\ s', \text{ and} \\ \text{If } th' \neq th &\text{ then } cprec\ s\ th' = cprec\ s'\ th'. \end{aligned}$$

This means again we do not have to recalculate the RAG and also not the current precedences for the other threads. Since th is not alive anymore in state s , there is no need to calculate its current precedence.

Set th prio: We assume that s' and $s \stackrel{def}{=} Set\ th\ prio::s'$ are both valid. We can show that

$$\begin{aligned} RAG\ s &= RAG\ s', \text{ and} \\ \text{If } th' \neq th \text{ and } th \notin \text{dependants } s\ th' &\text{ then } cprec\ s\ th' = cprec\ s'\ th'. \end{aligned}$$

The first property is again telling us we do not need to change the RAG . The second however states that only threads that are *not* dependants of th have their current precedence unchanged. For the others we have to recalculate the current precedence. To do

this we can start from th and follow the RAG -chains to recompute the $cprec$ of every thread encountered on the way using Lemma 2. Since the RAG is loop free, this procedure will always stop. The following two lemmas show, however, that this procedure can actually stop often earlier without having to consider all dependants.

*If $th \in \text{dependants } s \ th''$ and $cprec \ s \ th = cprec \ s' \ th$ then $cprec \ s \ th'' = cprec \ s' \ th''$.
If $th \in \text{dependants } s \ th'$, $th' \in \text{dependants } s \ th''$ and $cprec \ s \ th' = cprec \ s' \ th'$
then $cprec \ s \ th'' = cprec \ s' \ th''$.*

The first states that if the current precedence of th is unchanged, then the procedure can stop immediately (all dependent threads have their $cprec$ -value unchanged). The second states that if an intermediate $cprec$ -value does not change, then the procedure can also stop, because none of its dependent threads will have their current precedence changed.

$V \ th \ cs$: We assume that s' and $s \stackrel{def}{=} V \ th \ cs::s'$ are both valid. We have to consider two subcases: one where there is a thread to “take over” the released resource cs , and one where there is not. Let us consider them in turn. Suppose in state s , the thread th' takes over resource cs from thread th . We can show

$$RAG \ s = RAG \ s' - \{(C \ cs, T \ th), (T \ th', C \ cs)\} \cup \{(C \ cs, T \ th')\}$$

which shows how the RAG needs to be changed. This also suggests how the current precedences need to be recalculated. For threads that are not th and th' nothing needs to be changed, since we can show

$$\text{If } th'' \neq th \text{ and } th'' \neq th' \text{ then } cprec \ s \ th'' = cprec \ s' \ th''.$$

For th and th' we need to use Lemma 2 to recalculate their current precedence since their children have changed.

In the other case where there is no thread that takes over cs , we can show how to recalculate the RAG and also show that no current precedence needs to be recalculated.

$$\begin{aligned} RAG \ s &= RAG \ s' - \{(C \ cs, T \ th)\} \\ cprec \ s \ th' &= cprec \ s' \ th' \end{aligned}$$

$P \ th \ cs$: We assume that s' and $s \stackrel{def}{=} P \ th \ cs::s'$ are both valid. We again have to analyse two subcases, namely the one where cs is locked, and where it is not. We treat the second case first by showing that

$$\begin{aligned} RAG \ s &= RAG \ s' \cup \{(C \ cs, T \ th)\} \\ cprec \ s \ th' &= cprec \ s' \ th' \end{aligned}$$

This means we do not need to add a holding edge to the RAG and no current precedence needs to be recalculated.

In the second case we know that resource cs is locked. We can show that

$$RAG\ s = RAG\ s' \cup \{(T\ th, C\ cs)\}$$

$$\text{If } th \notin \text{dependants } s\ th' \text{ then } cprec\ s\ th' = cprec\ s'\ th'.$$

That means we have to add a waiting edge to the *RAG*. Furthermore the current precedence for all threads that are not dependants of *th* are unchanged. For the others we need to follow the edges in the *RAG* and recompute the *cprec*. However, like in the @case of text Set, this operation can stop often earlier, namely when intermediate values do not change.

A pleasing result of our formalisation is that the properties in this section closely inform an implementation of PIP: Whether the *RAG* needs to be reconfigured or current precedences need to be recalculated for an event is given by a lemma we proved.

6 Conclusion

The Priority Inheritance Protocol (PIP) is a classic textbook algorithm used in real-time operating systems in order to avoid the problem of Priority Inversion. Although classic and widely used, PIP does have its faults: for example it does not prevent deadlocks in cases where threads have circular lock dependencies.

We had two goals in mind with our formalisation of PIP: One is to make the notions in the correctness proof by Sha et al. [7] precise so that they can be processed by a theorem prover. The reason is that a mechanically checked proof avoids the flaws that crept into their informal reasoning. We achieved this goal: The correctness of PIP now only hinges on the assumptions behind our formal model. The reasoning, which is sometimes quite intricate and tedious, has been checked beyond any reasonable doubt by Isabelle/HOL. We can also confirm that Paulson’s inductive method for protocol verification [5] is quite suitable for our formal model and proof. The traditional application area of this method is security protocols. The only other application of Paulson’s method we know of outside this area is [10].

The second goal of our formalisation is to provide a specification for actually implementing PIP. Textbooks, for example [9, Section 5.6.5], explain how to use various implementations of PIP and abstractly discuss their properties, but surprisingly lack most details for a programmer who wants to implement PIP. That this is an issue in practice is illustrated by the email from Baker we cited in the Introduction. We achieved also this goal: The formalisation gives the first author enough data to enable his undergraduate students to implement PIP (as part of their OS course) on top of PINTOS, a small operating system for teaching purposes. A byproduct of our formalisation effort is that nearly all design choices for the PIP scheduler are backed up with a proved lemma. We were also able to establish the property that the choice of the next thread which takes over a lock is irrelevant for the correctness of PIP. Earlier model checking approaches which verified implementations of PIP [2,3,11] cannot provide this kind of “deep understanding” about the principles behind PIP and its correctness.

PIP is a scheduling algorithm for single-processor systems. We are now living in a multi-processor world. So the question naturally arises whether PIP has any relevance in such a world beyond teaching. Priority Inversion certainly occurs also in multi-processor systems. However, the surprising answer, according to [8], is that except for

one unsatisfactory proposal nobody has a good idea for how PIP should be modified to work correctly on multi-processor systems. The difficulties become clear when considering that locking and releasing a resource always requires a small amount of time. If processes work independently, then a low priority process can “steal” in such an unguarded moment a lock for a resource that was supposed allow a high-priority process to run next. Thus the problem of Priority Inversion is not really prevented. It seems difficult to design a PIP-algorithm with a meaningful correctness property on a multi-processor systems where processes work independently. We can imagine PIP to be of use in situations where processes are *not* independent, but coordinated via a master process that distributes work over some slave processes. However, a formal investigation of this is beyond the scope of this paper. We are not aware of any proofs in this area, not even informal ones.

The most closely related work to ours is the formal verification in PVS for Priority Ceiling done by Dutertre [1]. His formalisation consists of 407 lemmas and 2500 lines of “specification” (we do not know whether this includes also code for proofs). Our formalisation consists of around 210 lemmas and overall 6950 lines of readable Isabelle/Isar code with a few apply-scripts interspersed. The formal model of PIP is 385 lines long; the formal correctness proof 3800 lines. Some auxiliary definitions and proofs took 770 lines of code. The properties relevant for an implementation took 2000 lines. Our code can be downloaded from ...

References

1. B. Dutertre. The Priority Ceiling Protocol: Formalization and analysis using PVS. Technical report, System Design Laboratory, SRI International, Menlo Park, CA, October 1999. Available at <http://www.sdl.sri.com/dsa/publis/prio-ceiling.html>.
2. J. M. S. Faria. *Formal Development of Solutions for Real-Time Operating Systems with TLA+/TLC*. PhD thesis, University of Porto, 2008.
3. E. Jahier, B. Halbwegs, and P. Raymond. Synchronous Modeling and Validation of Priority Inheritance Schedulers. In *Proc. of the 12th International Conference on Fundamental Approaches to Software Engineering (FASE)*, volume 5503 of *LNCS*, pages 140–154, 2009.
4. B. W. Lampson and D. D. Redell. Experiences with Processes and Monitors in Mesa. *Communications of the ACM*, 23(2):105–117, 1980.
5. L. C. Paulson. The Inductive Approach to Verifying Cryptographic Protocols. *Journal of Computer Security*, 6(1–2):85–128, 1998.
6. G. E. Reeves. Re: What Really Happened on Mars? *Risks Forum*, 19(54), 1998.
7. L. Sha, R. Rajkumar, and J. P. Lehoczky. Priority Inheritance Protocols: An Approach to Real-Time Synchronization. *IEEE Transactions on Computers*, 39(9):1175–1185, 1990.
8. U. Steinberg, A. Bötcher, and B. Kauer. Timeslice Donation in Component-Based Systems. In *Proc. of the 6th International Workshop on Operating Systems Platforms for Embedded Real-Time Applications (OSPERT)*, pages 16–23, 2010.
9. U. Vahalia. *UNIX Internals: The New Frontiers*. Prentice-Hall, 1996.
10. J. Wang, H. Yang, and X. Zhang. Liveness Reasoning with Isabelle/HOL. In *Proc. of the 22nd International Conference on Theorem Proving in Higher Order Logics (TPHOLs)*, volume 5674 of *LNCS*, pages 485–499, 2009.
11. A. Wellings, A. Burns, O. M. Santos, and B. M. Brosgol. Integrating Priority Inheritance Algorithms in the Real-Time Specification for Java. In *Proc. of the 10th IEEE International*

- Symposium on Object and Component-Oriented Real-Time Distributed Computing (ISORC)*, pages 115–123. IEEE Computer Society, 2007.
12. V. Yodaiken. Against Priority Inheritance. Technical report, Finite State Machine Labs (FSMLabs), 2004.