

Final Report Option 1 – AOARD 104105

Formal System Verification for Trustworthy Embedded Systems

June Andronick, Gerwin Klein, Toby Murray

November 2011

Version 1 - September 2011 - submitted to AOARD Version 2 - November 2011 - includes final thesis version

1 Introduction

This report summarises the work done in Option 1 of AOARD project 104105, Formal System Verification for Trustworthy Embedded Systems. Option 1 of this project spans the time 16 Feb 2011 through 15 Sep 2011.

We begin by revisiting the original motivation, and the progress up to the start of Option 1. We continue by describing in more detail the outcomes of the Option 1 phase of the project in the form of a brief technical report. This part describes in particular the progress we have made in formalising a general framework that allows us to prove invariant properties about a system consisting of a microkernel, user-level trusted components, and user-level untrusted components. One of the key properties of this framework is that it will allow the integration of different formalisms for describing the behaviour of trusted components as well as the integration of other reasoning tools such as model checking into the verification.

We attach a publication describing the proof of the integrity security property of the seL4 microkernel [6]. The seL4 microkernel with its functional correctness proof [5] is the basis of this work. As we shall see in the technical part, the integrity property is one of the essential building blocks for composing separate component verifications into a whole-system proof. The seL4 microkernel is currently the only general-purpose OS kernel with a formal, machine-checked implementation-level proof of integrity.

We further attach the (to be submitted) draft of a BSc honours thesis [2] that demonstrates how a monad-based formalism can be used for the specification of functional behaviour of trusted components. This specification can interact and integrate with the specification of the seL4 kernel at multiple levels of abstraction. We have previously shown how a specification at a very high level of abstraction can be used to reason about the behaviour of trusted components [1]. We have also demonstrated in the seL4 proof that the same can be achieved down to the level of C code. We have yet to demonstrate in Option 2 of the project that refinement proofs can be conducted between these layers of abstraction, if possible independently of the abstraction layers of the kernel, and how they can be integrated into one overall theorem. The technical discussion in the rest of this report lays out how we plan to achieve this.

2 Background

As mentioned above the aim of our overall research agenda is to prove security and safety properties across a whole software system, including OS kernel, trusted, and untrusted components. Figure 1 shows a representation of this setup.

In previous work and previous parts of this project, we have developed a functional correctness proof for the seL4 microkernel, a prototype system demonstrating its use for dividing complex systems into small trusted and large untrusted components, a formalisation of that prototype system, and a security proof of it assuming properties of the trusted component and additional security properties of the seL4 kernel.

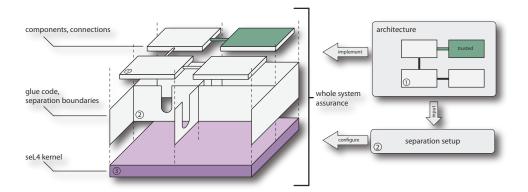


Figure 1: seL4-based system with trusted and untrusted components.

The aim of this AOARD project 104105 is to transport the security property to the implementation level of the overall system. Namely the goal is to

"investigate the creation of a formal overall system framework that enables the use of kernel correctness, component correctness and security analysis to conclude one formal, machine-checked theorem about the overall system on the level of its implementation."

Specifically, the objective is to

"investigate how to integrate a number of heterogeneous formal analyses into a single framework such that the overall effort for verification of the whole system is reduced to a provably correct combination of automatic, manual, and semi automatic partial proofs."

During Option 1 we have made significant progress on both goals by defining a formal automaton that describes the execution of the whole software system. It integrates kernel transitions with known behaviour transitions of trusted components and unknown behaviour transitions of untrusted components. Untrusted component behaviour is constrained by hardware mechanisms, mainly the virtual memory subsystem, and microkernel access control. The attached draft thesis report [2] shows an example instantiation of this automaton. Since automaton transitions in the framework are kept as general nondeterministic relations they can be instantiated by a number of different formalisms that admit different kinds of formal analysis. The default analysis method in this framework would be interactive proof in the theorem prover Isabelle/HOL. Integration of automated provers such as first-order provers and SMT-solvers is achieved by using the existing integration of such tools with the Isabelle theorem prover (with proof re-checking in Isabelle) or by using the external tool as a trusted oracle to Isabelle (without proof re-checking). Very successful examples of such integration with Isabelle and first-order automated provers or SMT solvers can now be found in the literature [3, 4]. If a particular abstraction level and its properties manage to

stay completely within the logic of any of these automatic tools, the corresponding proof can be automatic.

The technical part of the report below lays out how we plan to put the different parts of this framework together to obtain invariant properties of system execution, reducing the reasoning necessary to reasoning about system composition, kernel access control, and trusted component behaviour.

3 Whole-System Verification

Big picture. The aim is to be able to formally prove strong safety or security properties for large, complex software systems, typically comprising millions of lines of code. Our approach is to avoid verifying all of the code by designing the system as a set of components such that the targeted property only relies on the behaviour of a small number of components, so-called *trusted components*.

In other terms, trusted components can potentially violate the property and we need to prove that they do not—whereas untrusted components are not given enough rights to violate the property, therefore they can be as large and complex as needed and their behaviour can be arbitrary. The isolation or control of communications between components is ensured by the underlying kernel, seL4.

The kernel is the only component running in the privileged mode of the hardware, known as *kernel mode*. In this mode, the hardware allows free access to all resources. The kernel is therefore the most critical part of the system. Every other component runs in *user mode* and needs to perform a kernel call to be able to communicate with other components. User mode execution essentially only admits free access to the currently mapped memory and registers. All other interaction can be restricted via kernel access control. Which memory is currently mapped for each user component is again controlled by the kernel and thereby subject to the security policy of the system. If a user component tries to access memory it is not authorised to access, the hardware generates a page fault which leads to a kernel call. The kernel can then react appropriately. If the user component would like to invoke any other form of communication such as interrupts, synchronous or asynchronous message passing, it needs to invoke the corresponding kernel primitive. These are again subject to the access control policy.

This strictly enforced separation by policy allows us to reason about the execution of untrusted components without a specification of their behaviour.

Limitations. The line of reasoning presented above does not include timing or other side channels of the machine that user-level components may try to use for communication. By definition, side channels are those channels that are not visible in the formal model of the system. In our case, our highest-fidelity model includes the full functional behaviour of the kernel at the C-code level as well as the memory, register, and virtual memory subsystem behaviour of the machine. It does not include caches, the TLB, or the timing behaviour of the machine. It currently also does not yet include external devices. Any of these could potentially be used as communication side channels. How far that is possible depends mainly on the hardware and the implementation of the kernel. For these we still need to rely on traditional security analyses.

In relation to devices, the model presented below does not yet take into account interrupts apart from timer interrupts. We expect to be able to handle such device interrupts within the same framework in future work.

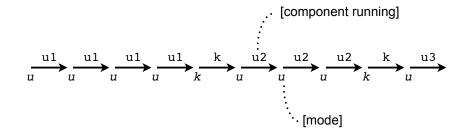
Assumptions. Some of the reasoning in the technical section below will lead to assumptions and thereby restrictions on the system policy. That means not all policies are valid for secure systems. This is to be expected: for instance, the policy that just give everyone access to everything is not a secure policy. The main notable assumptions are that we do not permit memory sharing between security domains and make restrictions on which communication primitives are permitted between domains. In future work we are looking to relax this assumption and allow limited forms of shared memory communication. A first formal set of assumptions is detailed in our paper on the integrity proof [6]. We expect further policy restrictions to arise from our work on confidentiality and system composition. Different levels of abstraction may also each bring their own restrictions with them.

Status. The status of the work described below is an initial semi-formal model and plan for formalising the rest of the framework. It is neither a fully formal model nor supported by formal proof yet. Instead it presents on a relatively detailed level our plan for achieving such a formalisation and our expectations on which proof obligations will remain for a full system analysis combining different levels of abstraction.

3.1 States, Transitions, System Execution

As mentioned, our aim is to prove that a given property holds at the source code level of a componentised system. Such a proof will generally be derived from a more general invariant about the whole system, i.e. a property I which needs to hold at each single step of any possible execution of the system. The idea is to pick an invariant that is strong enough to imply the desired security or safety property.

A typical execution of the system is represented, in a simplified version, in Figure 2: a user u_1 runs, then it performs a kernel call, which traps into kernel mode where the kernel executes the command, and then schedules another user, say u_2 ; this user then starts running until it in turn performs a kernel call, etc. The sequence depicted in Figure 2 is an execution trace generated by the system automaton containing the formalisation of kernel and user-component behaviour. The goal is to show properties over all possible execution traces of a given system. Each step in Figure 2 indicates which component is running, and each intermediate state indicates the current mode (user mode or kernel mode). For more precision, we will add two more pieces of information into the state. Firstly, for the beginning of each kernel execution, we note on whose behalf the kernel is running. In terms of the kernel, this is represented by the currently running thread. We do this





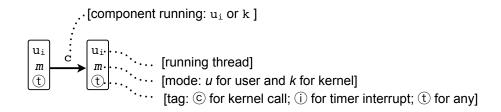


Figure 3: Representation of a system transition at the code level

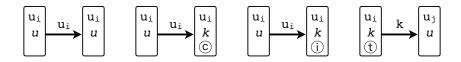


Figure 4: All possible system transitions at the code level

to distinguish kernel calls of trusted components from kernel calls of untrusted components. Secondly, a kernel execution is an atomic step, as it can't be interrupted. However, a user execution can be interrupted at any time when a timer interrupt occurs. This would result in a user transition ending in a state in kernel mode, followed by a kernel transition. To represent this, we will add a *tag* indicating if the switch to kernel mode is due to a call or a timer interrupt. We currently do not yet treat other kinds of interrupts on this level, such as device interrupts, although they are already modelled in the kernel correctness proofs. Figure 3 shows our representation of states and transitions, whereas Figure 4 shows all possible transitions in any seL4-based system. Using this more detailed representation of states, the example execution given in Figure 2 would be represented as in Figure 5. Now the example where an interrupt occurs, say during the execution of user u_1 , is represented in Figure 6.

So far we have illustrated the notion of system execution by example. We will keep using examples to illustrate the different theorems that will be used, but we will also try to show the

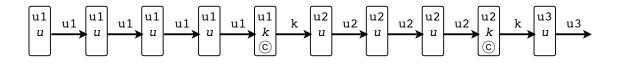


Figure 5: Example of system execution at the code level (detailed)

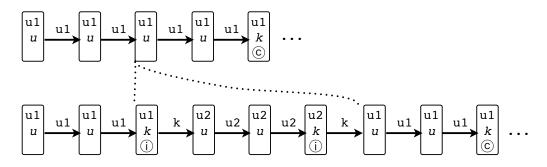


Figure 6: Example of system execution at the code level, including a timer interrupt

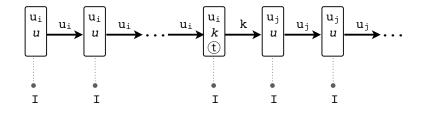


Figure 7: Representation of the general definition of a system execution, together with the property to prove, i.e. that the invariant I holds at each step of the execution.

general case as far as possible. In the case of execution, given the possible system transitions (Figure 4), the general definition would be as in Figure 7, i.e. an execution is a sequence of user transitions, followed by a kernel transition due to a kernel call or a time interrupt, followed by a new sequence of user transition for newly scheduled user, and so on and so forth. Figure 7 also illustrates the property we are aiming for, i.e. proving that the invariant I holds at each state of the execution.

3.2 Verification approach

The general idea for verifying that the invariant I holds at each step of the execution of the source code will be to, whenever possible, abstract the formal model of the execution at a high level where the verification will require less effort. The two main challenges in doing this will be: (1)

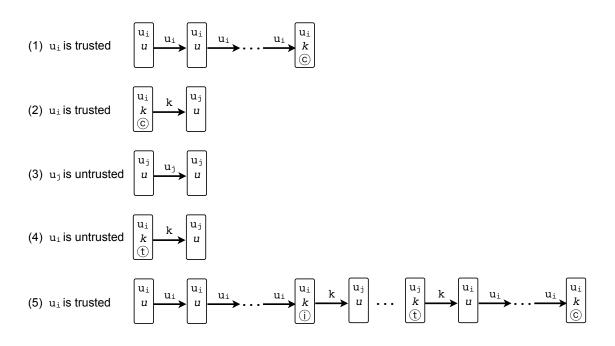


Figure 8: Decomposition of execution in 5 different cases

ensure that the proof at the high level is preserved by refinement, i.e. that I holds for any concrete level implementing the abstract high level; and (2) combine possibly different levels of abstraction and different formalisms used for various parts of the execution. Indeed, a different approach will be used to deal with different chunks of the execution; for instance, reasoning about trusted components execution will require more effort than reasoning about untrusted ones.

The plan is to decompose the possible executions described in Figure 7 into the 5 following cases, illustrated in Figure 8.

- 1. Trusted component execution up to a kernel call, with no interrupt.
- 2. Kernel call on behalf of a trusted component.
- 3. Single untrusted component transition.
- 4. Kernel call on behalf of an untrusted component or timer interrupt during execution of an untrusted component.
- 5. Trusted component execution up to a kernel call, but interleaved with interrupts leading to execution of other components.

The verification approach is that for each of these 5 cases, we will prove that the invariant holds after each state of the execution, assuming that it held when the execution started. The approach

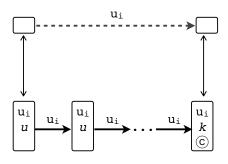


Figure 9: Refinement proof for trusted components

for this proof and its level of difficulty will be different for each case. Each one is described in the following.

3.3 Trusted component execution up to a kernel call, without interrupt (case 1)

Since, by definition, the trusted components have enough rights to potentially violate the targeted property, we need to model their behaviour precisely and prove that the property is not violated by this behaviour. This proof can be complex and ideally we will want to conduct it at a high level of abstraction. Therefore we need to (1) model an abstraction of the behaviour of the trusted components and prove that their implementation is a refinement of this model (illustrated in Figure 9); (2) prove that the abstract version I' of the invariant I holds for the abstract model (illustrated in Figure 10); and (3) prove that if the abstracted invariant I' holds on an abstract state, then I holds on its refined state (illustrated in Figure 11). Putting everything together, we will end up in a proof that the invariant I holds at the concrete level for each state of the execution of a trusted component up to a kernel call, as illustrated in Figure 12.

We have demonstrated refinement proofs down to the source-code level before in the verification of the seL4 microkernel. An additional challenge for trusted component verification is the step atomicity. While the kernel can safely be modelled as one atomic step, because interrupts are disabled, the execution of a trusted component could be interrupted at any time and we need to show that our invariant I does not only hold at begin and end of the possibly longer abstract execution, but at the end of each atomic step on the concrete level. This means the refinement proof obligation needs to be stronger than in our previous kernel verification.

Note that during the user-level execution of trusted components, only small parts of the invariant are expected to be relevant—those that talk about the state (memory and registers) of the running trusted component. The user-level execution of a trusted component will not be able to affect any other global system state yet. For such a change to occur, the component needs to perform a kernel call which leads us into the next case.

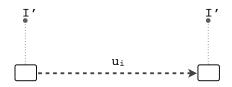


Figure 10: Invariant preservation proof at the abstract level

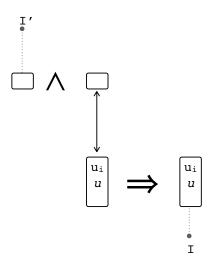


Figure 11: Proof that the invariant is preserved by refinement

3.4 Kernel call on behalf of a trusted component (case 2)

In this case, a trusted component has previously performed a computation and is about to place a kernel call. In terms of the system policy the trusted component will have enough authority for the kernel to permit actions in this call that may violate the security property and therefore the system invariant *I*. It is our proof obligation to show that this does not happen.

We noted in case 1 (trusted component execution) that only small parts of the invariant, if at all, will talk about the private state of the trusted component and therefore invariant preservation will not be too hard to prove. In case 2, however, the kernel call may have global effects. This is the interesting case in the invariant preservation proof. To show that these effects preserve the abstract invariant I', we need to use the result of the verification in case 1 to determine which parameters the kernel was invoked with. Given these parameters and the specification of the kernel at the suitable level of abstraction, we must show that I' is preserved. Per the existing kernel refinement theorem, we then get that I again holds after the kernel call has finished on the concrete level.

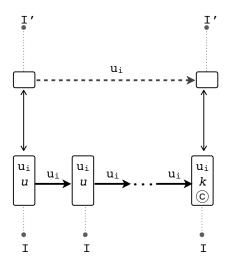


Figure 12: Proof scheme for case 1

3.5 Single untrusted component transition (case 3)

Untrusted components can potentially be large and complex. The idea is not to model any details of their behaviour at all, and to just use the fact that we have restricted their authority so they cannot violate the property. The authority each component has is described as a policy graph \mathcal{P} . This graph describes what authority each components has to other components or to memory regions. It is an abstraction of the actual detailed protection state given by the capabilities in the system. The policy graph also indicates which is the currently running component (called *subject*). When an untrusted component, say u_j , executes a *user* transition, the only thing that it can potentially modify or read is the memory it has access to. Any other actions, such as communicating with other components, require a kernel call, which is covered by the next case. Therefore the approach here is to prove, at an abstract level, that I' does not depend on the private state of the untrusted component, in particular its memory and registers.

We do this by first proving that if s'_1 is the (abstract) state before the untrusted user transition, and s'_2 the one after, then the only thing can possible have changed in s'_2 is u_j 's memory. This is noted $s'_2 =_{|u_j \dots mem} s'_1$, and is illustrated by Figure 13. This fact can be proved generically and once and for all.

Then we prove that I' does not depend on u_j_mem by showing that if it holds for any given state s'_1 , then it will hold on any state s'_2 that differs only on u_j_mem (as shown in Figure 14).

Now we prove that the abstract transition (modelled as *any* transition) is trivially refined by any implementation, including the actual source code of the untrusted component, as illustrated in Figure 15).

Putting all the theorems together, and using the invariant preservation by refinement (Fig-

$$\left(\underbrace{\mathbf{s}_{1'}}_{(\operatorname{any})} \underbrace{\mathbf{u}_{j}}_{\operatorname{s}_{2'}} \wedge \mathbf{u}_{j} = \operatorname{subject}(P(\mathbf{s}_{1'})) \wedge \mathbf{u}_{j}\operatorname{mem} = \operatorname{memory}(\mathbf{u}_{j}, P(\mathbf{s}_{1'})) \right)$$
$$\implies \mathbf{s}_{2'} = |u_{j}\operatorname{mem} \mathbf{s}_{1'}$$

Figure 13: Resulting state for any untrusted user transition

$$\left(s_{2'} =_{|uj_mem|} s_{1'} \land I'(s_{1'}) \right) \implies I'(s_{2'})$$

Figure 14: Invariant does not depend on untrusted user memory

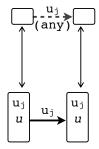


Figure 15: Refinement proof for untrusted components

ure 11), we obtain that the invariant I is trivially preserved by untrusted user transitions, as illustrated in Figure 16).

3.6 Kernel call or interrupt when untrusted component is running (case 4)

The same way we proved that the invariant does not depend on untrusted component executions, we need to prove that it neither depends on the kernel actions performed when those components invoke a kernel call. This case is more complicated than a user transition, because the kernel runs in privileged mode so it could potentially do more than what the running thread is authorised to

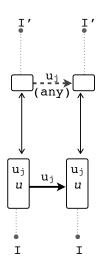


Figure 16: Proof scheme for case 3

$$\left(\underbrace{[\mathbf{s}_{1'}] \cdots \mathbf{s}_{\mathbf{s}_{2'}}}_{\mathbf{s}_{2'}} \wedge \mathbf{u}_{j} = \operatorname{subject}(P(\mathbf{s}_{1'})) \wedge \mathbf{u}_{j} \mathbb{W}_{auth} = \operatorname{authorized}_{write}(\mathbf{u}_{j}, P(\mathbf{s}_{1'})) \right)$$
$$\implies \mathbf{s}_{2'} = |u_{j} \mathbb{W}_{auth} \mathbf{s}_{1'}$$

Figure 17: Integrity theorem

do. A crucial property is therefore that it does not. The second crucial property is that the kernel call does not change the security policy. Again, the kernel potentially has the ability to do this and the proof obligation is that it does not. Both properties can be proved once and for all with the policy as a parameter. Our proof of integrity [6] does precisely this. It guarantees that the kernel will not modify objects that the running thread is not allowed to modify. In other words, if s'_1 is the (abstract) state before the kernel call — on behalf of an untrusted user u_j — and s'_2 the one after, then the only thing that is allowed to have changed in s'_2 is what u_j has write authority to as defined by the policy \mathcal{P} . This is noted $s'_2 =_{|u_j,W,auth} s'_1$, and is illustrated by Figure 17.

We then need to prove that the (abstracted) invariant I' does not depend on what the untrusted components are allowed to modify, as illustrated in Figure 18.

The invariant I' is likely to mention the policy and/or the protection state of the system. We therefore need the second property mentioned above, *authority confinement*. In particular, authority confinement implies that an untrusted component cannot change the policy graph, and therefore the

$$\left(\mathbf{s}_{2'} =_{|uj_W_auth} \mathbf{s}_{1'} \land \mathbf{I}'(\mathbf{s}_{1'})\right) \implies \mathbf{I}'(\mathbf{s}_{2'})$$

Figure 18: Invariant does not depend on what untrusted users can modify

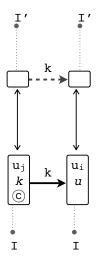


Figure 19: Proof scheme for case 4

policy graph is unchanged through all untrusted component transitions. This enables the theorems on such transition to be chained into longer executions.

Again, we now need to prove that the abstract kernel transition is refined by the kernel source code transition and therefore the concrete level preserves I. This is given by our previous seL4 correctness proof [5]. Putting all the theorems together, we obtain that the invariant I is preserved by kernel calls on behalf of untrusted user transitions, as illustrated in Figure 19.

The only proof obligation in this case that has to be proved per system and that is not yet discharged once and for all is that the invariant does not depend on changes untrusted components are authorised to perform. If the invariant is written appropriately, this should be a low-effort, potentially automatic proof.

The case where the switch to kernel mode is due to a timer interrupt instead of a kernel call follows precisely the same reasoning and should be even simpler, because in a timer interrupt, the kernel does not perform any authority-relevant actions.

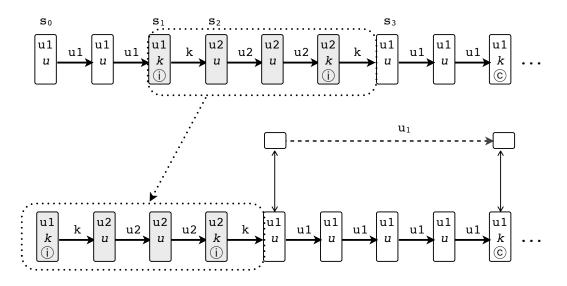


Figure 20: Reordering theorem, for case 5

3.7 Trusted component execution up to a kernel call, with interrupts (case 5)

In the last case an interrupt occurs during the execution of a trusted component, say u_1 , leading to the execution of another user u_2 being interleaved into the execution of u_1 (as we have seen in Figure 6). As explained in Section 3.3, our approach to reason about trusted components execution was to do the verification of trusted components behaviour at an abstract level, where it requires less effort (as illustrated in Figure 12). The problem is that a sequence of steps in the concrete level corresponds to a single bigger step in the abstract level. Therefore, we can only keep the reasoning at the abstract level if we observe this exact sequence of transitions at the low level. This will not be the case anymore if a timer interrupt occurs during this execution.

The approach here will thus be to prove that we can equivalently *re-order* the execution, so that we keep all u_1 steps in sequence to allow our refinement to hold. Namely, we want to prove that the invariant *I* holds at each step of the execution if and only if it holds on the execution where the interrupt occurred earlier, as illustrated in Figure 20. The argument is that: (1) timer interrupts do not change any memory content, and (2) u_1 transitions only change u_1 memory, which is disjoint from u_2 memory, because we assume there is no memory sharing between components. Therefore, the u_2 memory is the same in state s_2 as in state s_1 and again the same as in s_0 ; therefore the behaviour of u_2 should be the same whether it starts executing in s_2 or s_0 . Conversely, the execution of u_2 can only affect u_2 private memory and u_1 will therefore find the same state as in its execution before the reordering.

The key observation to note is that the synchronisation points in the system are the kernel calls of components. Re-ordering over timer interrupts is easy, as long as we manage to prove that the kernel protects integrity, authority confinement and confidentiality during this event. The first two we have proved, the latter is in progress. Re-ordering over kernel calls of trusted components is not always possible in general—the trusted component may set the system policy into a new state, for instance removing the whole untrusted component u_2 from the system. This fact would clearly be observable and the re-ordered execution would not be equivalent to the original one. We therefore plan to re-order the interleaved execution of u_2 before the execution of u_1 . This takes care of u_1 kernel calls.

The most interesting case in the reordering occurs when u_2 performs a kernel call at the end of its execution (instead of another timer interrupt occurring). In this case we need to argue that this kernel call does not affect the private execution of u_1 . Again we argue from the policy \mathcal{P} through the integrity theorem and authority confinement that u_2 does not possess enough authority to do anything that affects the private memory or registers of u_1 . The kernel call does potentially change kernel state that may affect the next kernel call of the next component, but our re-ordering preserves the order of kernel calls—it just moves u_1 executions over u_2 executions or u_2 kernel calls to the right. By not changing the order of kernel calls our re-ordering preserves the order of synchronisation events in the system and makes sure that each kernel call sees equivalent initial states before and after re-ordering.

The same argument applies if u_2 is itself another trusted component. For this we additionally need the proof obligation that trusted components do not make system policy changes that are observable in the private user-level execution of another trusted component.

Equivalently re-ordering executions this way enables us to deal with the first part of the (reordered) execution using case 3 for u_2 (and cases 2 and 4 for the kernel transitions) and deal with the second part for u_1 using case 1.

4 Conclusion

In this report we have summarised the work of Option 1 in AOARD project 104105 from 16 Feb 2011 through 15 Oct 2011. We have made significant progress on a formal framework for the verification of entire microkernel-based software systems. We have outlined the formal framework in the technical part of the report, and we have attached a project publication on two fundamental kernel properties that enable the use of this framework. The attached thesis draft shows how a user-level specification can integrate with kernel specifications at different levels of abstraction.

References

 June Andronick, David Greenaway, and Kevin Elphinstone. Towards proving security in the presence of large untrusted components. In Gerwin Klein, Ralf Huuck, and Bastian Schlich, editors, 5th SSV, Vancouver, Canada, Oct 2010. USENIX.

- [2] Nelson Billing. Formal functional specification of a security critical component in Isabelle/HOL. BE thesis, School Comp. Sci. & Engin., University NSW, Sydney 2052, Australia, Nov 2011. Draft, to be submitted.
- [3] Jasmin Christian Blanchette, Sascha Böhme, and Lawrence C. Paulson. Extending sledgehammer with SMT solvers. In Nikolaj Bjørner and Viorica Sofronie-Stokkermans, editors, *CADE*, volume 6803 of *Lecture Notes in Computer Science*, pages 116–130. Springer, 2011.
- [4] Sascha Böhme and Tobias Nipkow. Sledgehammer: Judgement day. In Jürgen Giesl and Reiner Hähnle, editors, *IJCAR*, volume 6173 of *Lecture Notes in Computer Science*, pages 107–121. Springer, 2010.
- [5] Gerwin Klein, June Andronick, Kevin Elphinstone, Gernot Heiser, David Cock, Philip Derrin, Dhammika Elkaduwe, Kai Engelhardt, Rafal Kolanski, Michael Norrish, Thomas Sewell, Harvey Tuch, and Simon Winwood. seL4: Formal verification of an operating system kernel. *CACM*, 53(6):107–115, Jun 2010.
- [6] Thomas Sewell, Simon Winwood, Peter Gammie, Toby Murray, June Andronick, and Gerwin Klein. seL4 enforces integrity. In Marko C. J. D. van Eekelen, Herman Geuvers, Julien Schmaltz, and Freek Wiedijk, editors, 2nd ITP, volume 6898 of LNCS, pages 325–340, Nijmegen, The Netherlands, Aug 2011. Springer-Verlag.

Attachments

seL4 Enforces Integrity

Thomas Sewell¹, Simon Winwood^{1,2}, Peter Gammie¹, Toby Murray^{1,2}, June Andronick^{1,2}, and Gerwin Klein^{1,2}

¹ NICTA, Sydney, Australia^{*}

² School of Computer Science and Engineering, UNSW, Sydney, Australia

{first-name.last-name}@nicta.com.au

Abstract. We prove that the seL4 microkernel enforces two high-level access control properties: integrity and authority confinement. Integrity provides an upper bound on write operations. Authority confinement provides an upper bound on how authority may change. Apart from being a desirable security property in its own right, integrity can be used as a general framing property for the verification of user-level system composition. The proof is machine checked in Isabelle/HOL and the results hold via refinement for the C implementation of the kernel.

1 Introduction

Enforcing access control is one of the primary security functions of an operating system (OS) kernel. Access control is usually defined as two properties: confidentiality, which means that information is not acquired without read authority, and integrity, which means that information is not modified without write authority. These properties have been well studied in relation to classical security designs such as the Bell-LaPadula model [3]. For dynamic access control systems, such as the capability system in the seL4 microkernel, an additional property is of interest: authority confinement, which means that authority may not be spread from one subject to another without explicit transfer authority.

We have previously verified the functional correctness of seL4 [12]. In this work we prove that seL4 correctly enforces two high level security properties: integrity and authority confinement.

We define these properties with reference to a user-supplied security policy. This policy specifies the maximum authority a system component may have. Integrity limits state mutations to those which the policy permits the subject components to perform. Authority confinement limits authority changes to those where components gain no more authority than the policy permits. The policy provides mandatory access control bounds; within these bounds access control is discretionary.

While integrity is an important security property on its own, it is of special interest in formal system verification. It provides a framing condition for the

^{*} NICTA is funded by the Australian Government as represented by the Department of Broadband, Communications and the Digital Economy and the Australian Research Council through the ICT Centre of Excellence program

execution of user-level components, telling us which parts of the system do not change. In a rely-guarantee framework, the integrity property gives us useful guarantee conditions for components without the need to consult their code. This is because the kernel, not the component, is providing the guarantee. This becomes especially important if the system contains components that are otherwise beyond our means for formal code-level verification, such as a Linux guest operating system with millions of lines of code. We can now safely and formally compose such parts with the rest of the system.

Access control properties and framing conditions have been extensively studied. Proving these properties about a real OS kernel implementation, however, has not been achieved before [10]. Specifically, the novelty and contributions of this work are:

- The Isabelle/HOL formalisation and generalisation of integrity and authority confinement for a real microkernel implementation.
- To the best of our knowledge the first code-level proof of high-level access control properties of a high-performance OS kernel.

Our proof is connected to reality via refinement to the C implementation. This means we must deal with all the complexities and corner cases of the kernel we have, rather than laying out a kernel design which fits neatly with our desired access control model and hoping to implement it later.

We make one kind of simplifying assumption: we place restrictions on the policy. These forbid some kinds of interaction between components which are difficult for us to reason about. Although we have not yet applied the theorem to a large system, our choice of assumptions has been guided by a previous case study [2] on a secure network access device (SAC), with a dynamic and realistic security architecture.

We are confident that a significant variety of security designs will, after some cosmetic adjustments, comply with our restrictions. We support fine grained components, communication between them via memory sharing and message passing, delegation of authority to subsystems and dynamic creation and deletion of objects. We support but restrict propagation of authority and policy reconfiguration at runtime.

In the following, Sect. 2 gives a brief introduction to access control in general and to the seL4 access control system in particular. Sect. 2 also introduces a part of the aforementioned SAC system as a running example. Sect. 3 gives a summary of the formalisation as well as the final Isabelle/HOL theorems and Sect. 4 discusses the results together with our experience in proving them.

2 Access Control Enforcement and seL4

This section introduces relevant concepts from the theory of access control and our approach to instantiating them for seL4. For ease of explanation we will first introduce a running example, then use it to describe the seL4 security mechanisms available, and then compare to the theory of access control.

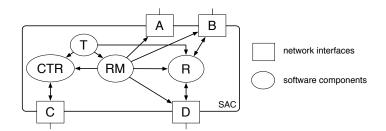


Fig. 1. System snapshot, routing between data D and back-end B network.

2.1 Example

The secure access controller (SAC) was designed in an earlier case study of ours [2] and will serve as a running example in this work. The purpose of the SAC is to switch one front-end terminal between different back-end networks of separate classification levels. The security goal of the system is to avoid information flow between the back-end networks.

Fig. 1 shows the main components of the SAC as nodes and their logical authority connections as edges. To the outside, the SAC provides four network interfaces: two back-end networks A and B, one control network C, and one data network D. Networks C and D are attached to a front-end terminal. The purpose of the SAC is to connect either A to D or B to D at a time without information flow between A and B. Internally, we have four components: a timer T, a controller user interface CTR, the router manager RM, and a router instance R. The router manager RM will upon a switch request, tear down R, remove all access from it, create and start a fresh R component, and connect it to the requested networks. RM is a small, trusted component and has access to all networks. Verification needs to show that it does not abuse this access. The router R, on the other hand, is a large, untrusted instance of Linux. It will only ever be given access to one of the back-end networks during its lifetime.

In previous work [2], we have shown that, assuming a specific policy setup, assuming correct behaviour of RM, and assuming that the kernel correctly enforces access control, the security goal of the system will be enforced.

The work in this paper helps us to discharge the latter two assumptions: we can use the integrity property as a framing condition for verifying the behaviour of RM, and we can use the same property to make sure that R stays within its information flow bounds. To complete the verification, we would additionally need the confidentially side of access control as well as a certified initial policy set up. Both are left for future work.

2.2 seL4

The seL4 microkernel is a small operating system kernel. As a microkernel, it provides a minimal number of services to applications: interprocess communication, threads, virtual memory, access control, and interrupt control. As mentioned, seL4 implements a capability-based access control system [6]. Services, provided by a set of methods on kernel implemented objects, are invoked by presenting to the kernel a capability that refers to the object in question and has sufficient access rights for the requested method. For the purposes of this paper, the following four object classes are the most relevant.

- **CNodes** Capabilities are stored in kernel-protected objects called *CNodes*. These CNodes can be composed into a *CSpace*, a set of linked CNodes, that defines the set of capabilities possessed by a single thread. CNode methods allow copying, insertion and removal of capabilities. For a thread to use a capability, this capability must be stored in the thread's CSpace. CNodes can be shared across CSpaces. The links in Fig. 1 mean that the collective CSpaces of a component provide enough capabilities to access or communicate with another component.
- Virtual Address Space Management A virtual address space in seL4 is called a *VSpace*. In a similar way to CSpaces, a VSpace is composed of objects provided by the microkernel. On ARM and Intel IA32 architectures, the root of a VSpace consists of a *Page Directory* object, which contains references to *Page Table* objects, which themselves contain references to *Frame* objects representing regions of physical memory. A Frame can appear in multiple VSpaces, and thereby implement shared memory between threads or devices such as the SAC networks.
- **Threads** Threads are the unit of execution in seL4, the *subjects* in access control terminology. Each thread has a corresponding TCB (thread control block), a kernel object that holds its data and provides the access point for controlling it. A TCB contains capabilities for the thread's CSpace and VSpace roots. Multiple threads can share the same CSpace and VSpace or parts thereof. A component in the SAC example may consist of one or more threads.
- **Inter-process Communication (IPC)** Message passing between threads is facilitated by *Endpoints* (EP). The kernel provides *Synchronous Endpoints* with rendezvous-style communication, and *Asynchronous Endpoints* (AEP) for notification messages. Synchronous endpoints can also be used to transfer capabilities if the sender's capability to the endpoint has the *Grant* right. The edges in the SAC example of Fig. 1 that are physical communication links, are implemented using endpoints.

The mechanisms summarised above are flexible, but low-level, as customary in microkernels. Fig. 2 shows parts of the implementation of the link between the CTR and RM components of Fig. 1. The link is implemented via a synchronous endpoint EP. The CTR and RM components both consist of one main thread, each with a CSpace containing a capability to the endpoint (among others), and each with a VSpace containing page directories (pd), page tables (pt), and frames f_n implementing private memory.

These mechanisms present a difficulty for access control due to their fine grained nature. Once larger components are running, there can easily exist hundreds of thousands of capabilities in the system. In addition, seL4 for ARM has

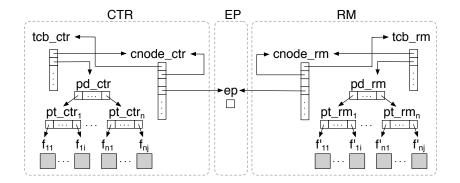


Fig. 2. SAC capabilities (partial).

15 different capability types, several of which have specific rights and variations. We will prune this complexity down by making some simplifying observations. One such observation is that many kinds of capabilities will not be shared without the sharing components trusting each other. For example sharing part of a CSpace with an untrusted partner makes little sense, as seL4 does not provide a mechanism for safely using a capability of unknown type.

2.3 Access Control Enforcement

An access control system controls the access of *subjects* to *objects* [14], by restricting the operations that subjects may perform on objects in each state of the system. As mentioned above, in seL4 the subjects are threads, the objects are all kernel objects, including memory pages and threads themselves.

The part of the system state used to make access control decisions, i.e., the part that is examined by the kernel to decide which methods each subject may perform on which object, is called the *protection state*. In seL4, this protection state is mostly represented explicitly in the capabilities present in the CSpace of each subject. This explicit, fine-grained representation is one of the features of capability-based access control mechanisms. In reality, however, some implicit protection state remains, for instance encoded in the control state of a thread, or in the presence of virtual memory mappings in a VSpace.

The protection state governs not only what operations are allowed to be performed, but also how each subject may modify the protection state. For instance, the authority for capabilities to be transmitted and shared between subjects is itself provided by CNode or endpoint capabilities.

As seen previously in Fig. 2, the protection state of a real microkernel can be very detailed, and therefore cumbersome to describe formally. It is even more cumbersome to describe precisely what the allowed effects of each operation are at this level. Hence we make use of the traditional concept of a policy which can be seen as an abstraction of the protection state: we assign a label to each object and subject, and we specify the authority between labels as a directed graph. The abstraction also maps the many kinds of access rights in seL4 protection states

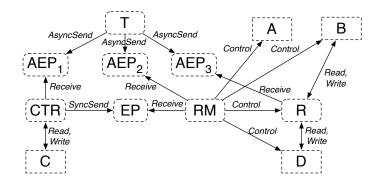


Fig. 3. SAC authority (except self-authority).

into a simple enumerated set of authority types. This simplifies the setup in three ways: the number of labels can be much smaller than the number of objects in the system, the policy is static over each system call whereas the protection state may change and, finally, we can formulate which state mutations are allowed by consulting the policy, rather than the more complex protection state.

The abstraction level of the policy is constrained only by the wellformedness assumptions we make in Sect. 3. Within these bounds it can be chosen freely and suitably for any given system.

Fig. 3 shows an abstract policy, only mildly simplified for presentation, that corresponds to a possible protection state of the SAC at runtime. The objects in the system are grouped by labels according to the intention of the component architecture. The RM label, for instance, includes all objects in the RM box of Fig. 2. The communication endpoints between components have their own label to make the direction of information flow explicit. The edges in the figure are annotated with authority types described in Sect. 3.

Correct access control enforcement, with respect to a security policy, can be decomposed into three properties about the kernel: *integrity, confidentiality* and *authority confinement*. The two former properties are the same as introduced in Sect. 1, only the notion of what is permitted is taken from the policy rather than the protection state. The latter property ensures the current subject cannot escalate its authority (or another subject's) above what the policy allows.

Note that we may have some components in the system, such as RM in the SAC, which have sufficient authority to break authority confinement. The authority confinement theorem will assume these components are not the current subject, and we will be obliged to provide some other validation of their actions.

3 Formalisation of Integrity Preservation

This section sketches our Isabelle/HOL formalisation of the integrity and authority confinement properties. While it is impossible in the space constraints of a paper to give the full detail, we show the major definitions and the top level theorems to provide a flavour of the formalisation. For formal details on the kernel-level concepts beyond our description below we refer the interested reader to the published Isabelle/HOL specification of the kernel [1] that the definitions here build on.

We have already introduced the notion of a policy, an upper bound on the protection state of the system, and an accompanying abstraction, a mapping from detailed subject and object names up to a smaller set of component labels. We roll these objects together with a subject label into the PAS record (policy, abstraction, subject) which is an input to all of our access control predicates.

The type parameter 'l here is any convenient type for component labels. The policy field is a graph (a set of triples 'l × auth × 'l), whose vertices are policy labels and whose edges are authority types from the type auth which will be discussed shortly. Abstraction functions are provided for seL4's namespaces: objects (i.e. system memory), interrupt request numbers and address space identifiers. Each of these is mapped up to a policy label.

The subject in the PAS record identifies the label of the current subject. We must associate all (write) actions with a subject in order to define the integrity property we are proving. We will pick the subject associated with any kernel actions at kernel entry time, choosing the label of the currently running thread. This coarse division of actions between subjects causes some problems for message transfers, as will be discussed below.

The identification of a subject makes all of our access control work *subjective*. The integrity property depends on which subject carries out an action, because whether that action is allowed or not depends on the allowed authority of the subject performing it. Wellformedness of PAS records, encapsulating our policy assumptions, will be defined in a subjective manner as well, so for any given system and policy the authority confinement proof may be valid only for less-trusted subjects which satisfy our policy assumptions.

3.1 Authority Types

We made the observation in Sect. 2 that most kinds of objects in seL4 are not shared by mutually distrusting components. We found that the objects that could be safely shared were those where authority to that object could be partitioned between capabilities. Endpoints are a good example: Capabilities to endpoints can be send-only or receive-only. Every endpoint in the SAC has a single sending component and a single receiving component, as is typical in seL4 system architectures. Memory frames may also be read only or read-write and a typical sharing arrangement has a single writer, though possibly many readers.

This led us to our chosen formalisation of authority types.

datatype auth = Receive | SyncSend | AsyncSend | Reset | Grant | Write | Read | Control

The Receive, Read and Write authorities have been described above. We distinguish endpoint types via SyncSend and AsyncSend. This distinction anticipates future work on confidentiality, where synchronous sends spread information in both directions. Capabilities to endpoints may also have the Grant right which permits other capabilities to be sent along with messages. The Reset authority is conferred by all capabilities to endpoints, since these capabilities can sometimes cause an endpoint reset even if they have zero rights. Finally, the Control authority is used to represent complete control over the target; it is a conservative approximation used for every other kind of authority in the system.

3.2 Subjective Policy Wellformedness

We aim to show authority is confined by a policy. This will not be true for all kinds of policies. If a policy gives the subject a **Grant** authority to any other component in addition to some authority which the other component should not have, that policy can clearly be invalidated. We define wellformedness criteria on policies, and thereby system architectures, given a specific subject, as follows. In our example Fig. 3, we would expect the policy to be wellformed for the subjects CTR, T, and R, but not RM.

 $\begin{array}{l} \mbox{policy-wellformed } policy \ irqs \ subject \equiv \\ (\forall a. \ (subject, \ {\sf Control}, \ a) \in policy \longrightarrow subject = a) \land \\ (\forall a. \ (subject, \ a, \ subject) \in policy) \land \\ (\forall s \ r \ ep. \\ (s, \ {\sf Grant}, \ ep) \in policy \land (r, \ {\sf Receive}, \ ep) \in policy \longrightarrow \\ (s, \ {\sf Control}, \ r) \in policy \land (r, \ {\sf Control}, \ s) \in policy) \land \\ (\forall i \in irqs. \ \forall p. \ (i, \ {\sf AsyncSend}, \ p) \in policy \longrightarrow (subject, \ {\sf AsyncSend}, \ p) \in policy) \end{array}$

The first requirement is that the subject cannot have **Control** authority over another component. If it did there would be no point in separating these components, as the subject might coerce the target into taking actions on its behalf.

The second requirement is that the subject has all kinds of authority to itself. We always consider components permitted to reconfigure themselves arbitrarily.

The Grant restriction observes that successful capability transfers over messages are as problematic for access control as Control authority. In each direction this restriction could be lifted if we introduced more complexity.

In the sending direction the problem is that the sender can transfer an arbitrary capability into the receiver's capability space, giving the sender the new authority to rescind capabilities from the receiver's capability space in the future. It may be possible in seL4 for a receiver to partition its capability space to make this safe, but we know of no use case that justifies the resulting complexity.

In the receiving direction the problem is in the way we fix the subject of the message send. Synchronous sends in seL4 complete when both sender and receiver are ready. If the sender is ready when our subject makes a receive system call, it may appear that the receiver has broken authority confinement by magically acquiring new authority. In fact the authority belonged to the sender, which was involved as a subject in some sense, but not in a manner that is easy to capture.

The final policy assumption relates to interrupts. Interrupts may arrive at any time, delivering an asynchronous message to a thread waiting for that interrupt. We must allow the current subject to send this message. We hope to revisit our simple notion of the current subject in future work.

3.3 Policy/Abstraction Refinement

We define a kernel state s as being a refinement of the policy p as follows.

```
\begin{array}{l} \mathsf{pas-refined} \ p \ s \equiv \\ \mathsf{policy-wellformed} \ (\mathsf{pasPolicy} \ p) \ (\mathsf{range} \ (\mathsf{pasIRQAbs} \ p)) \ (\mathsf{pasSubject} \ p) \ \land \\ \mathsf{irq-map-wellformed} \ p \ s \ \land \\ \mathsf{auth-graph-map} \ (\mathsf{pasObjectAbs} \ p) \ (\mathsf{state-objs-to-policy} \ s) \subseteq \mathsf{pasPolicy} \ p \ \land \\ \mathsf{state-asids-to-policy} \ p \ s \subseteq \mathsf{pasPolicy} \ p \ \land \\ \mathsf{state-irqs-to-policy} \ p \ s \subseteq \mathsf{pasPolicy} \ p \end{array}
```

The kernel state refines the policy if the various forms of authority contained within it, when labelled by the abstraction functions, are a subset of the policy. The full definitions of the extraction functions for authority from the kernel state are too detailed to describe here. In summary, a subject has authority over an object for one of these reasons:

- it possesses a capability to the object.
- it is a thread which is waiting to conclude a message send. For performance reasons the capability needed to start the send is not rechecked on completion, and thus the thread state is an authority in its own right.
- it possesses the parent capability in the capability derivation tree (cdt) of a capability stored in the object.
- the page tables link the subject to the object.
- the active virtual address space database names the object as the page directory for an address space identifier the subject owns.
- the interrupt despatch mechanism lists the object as the receiver for an interrupt request number the subject owns.

Note that none of this authority extraction is subjective. The **pas-refined** predicate is subjective only because it also asserts **policy-wellformed**. The reason for this is convenience: these properties are almost always needed together in the proof.

3.4 Specification of Integrity

We define access control integrity subjectively as follows:

integrity $p \ s \ s' \equiv$

- $(\forall x. \text{ object-integrity } p \text{ (pasObjectAbs } p x) \text{ (kheap } s x) \text{ (kheap } s' x)) \land$
- $(\forall x. memory-integrity p x (tcb-states-of-state s) (tcb-states-of-state s')$
 - (auth-ipc-buffers s) (memory-of s x) (memory-of s' x)) \land
- $(\forall x. \text{ cdt-integrity } p \ x \ (\text{cdt} \ s \ x, \text{ is-original-cap } s \ x) \ (\text{cdt} \ s' \ x, \text{ is-original-cap } s' \ x))$

This says that a transition from s to s' satisfies access control integrity if all kernel objects in the kernel object heap (kheap s), user memory and the capability derivation tree (cdt s) were changed in an acceptable manner.

The object level integrity predicate is defined by eight introduction rules. The following three rules are representative:

 $\begin{array}{ll} x = \mathsf{pasSubject} \ p & ko = ko' \\ \hline \mathsf{object-integrity} \ p \ x \ ko \ ko' & \mathsf{object-integrity} \ p \ x \ ko \ ko' \\ \hline \mathsf{bo} = \mathsf{Some} \ (\mathsf{TCB} \ tcb) & ko' = \mathsf{Some} \ (\mathsf{TCB} \ tcb') \\ \exists \ ctxt'. \ tcb' = \ tcb(|\mathsf{tcb-context} \ := \ ctxt', \ \mathsf{tcb-state} \ := \ \mathsf{Running}|) \\ \mathsf{receive-blocked-on} \ ep \ (\mathsf{tcb-state} \ tcb) & \mathsf{auth} \in \{\mathsf{SyncSend}, \ \mathsf{AsyncSend}\} \\ (\mathsf{pasSubject} \ p, \ \mathsf{auth}, \ \mathsf{pasObjectAbs} \ p \ ep) \in \mathsf{pasPolicy} \ p \end{array}$

object-integrity $p \ l' \ ko \ ko'$

These cases allow the subject to make any change to itself, to leave anything unchanged, and to send a message through an endpoint it has send access to and into the registers (called the tcb-context here) of a waiting receiver. Note that it is guaranteed that the receiver's registers are changed only if the message transfer completes and that the receiver's state is changed to Running. It is likewise guaranteed by memory-integrity that a receiver's in-memory message buffer is changed only if the message transfer completes, which is the reason for the complexity of the arguments of the memory-integrity predicate above.

The additional **object-integrity** cases include a broadly symmetric case for receiving a message and resuming the sender, for resetting an endpoint and evicting a waiting sender or receiver, for updating the list of threads waiting at an endpoint, and for removing a virtual address space the subject owns from the active virtual address space database.

The cdt-integrity predicate limits all cdt changes to the subject's label.

The crucial property about integrity is transitivity:

Lemma 1. integrity $p \ s_0 \ s_1 \land$ integrity $p \ s_1 \ s_2 \longrightarrow$ integrity $p \ s_0 \ s_2$

This must be true at the top level for our statement about a single system call to compose over an execution which is a sequence of such calls.

Integrity is also trivially reflexive:

Lemma 2. integrity $p \ s \ s$

3.5 Top Level Statements

Both integrity and pas-refined should be invariants of the system, which can be demonstrated using Hoare triples in a framework for reasoning about state monads in Isabelle/HOL. We have previously reported on this framework in depth [5]. In summary, the precondition of a Hoare triple in this framework is a predicate on the pre-state, the post condition is a predicate on the return value of the function and the post-state. In the case below, the return value is *unit* and can be ignored in the post condition. The framework provides a definition, logic, and automation for assembling such triples. **Theorem 1 (Integrity).** The property integrity pas st holds after all kernel calls, assuming that the protection state refines the policy, assuming the general system invariants invs and ct-active (current thread is active) for non-interrupt events, assuming the policy subject is the current thread, and assuming that the kernel state st is the state at the beginning of the kernel call. In Isabelle:

{pas-refined $pas \cap invs \cap (\lambda s. ev \neq \text{Interrupt} \longrightarrow \text{ct-active } s) \cap is-subject <math>pas \circ \text{cur-thread} \cap (\lambda s. s = st)$ } call-kernel ev{ λ -. integrity pas st}

We have shown in previous work [12] that the preconditions invs and $ev \neq$ Interrupt \longrightarrow ct-active s hold for any system execution at kernel entry.

Theorem 2 (Authority Confinement). The property pas-refined pas is invariant over kernel calls, assuming again the general system invariants invs and ct-active for non-interrupt events, and assuming that the current subject of the policy is the current thread.

{pas-refined $pas \cap invs \cap (\lambda s. ev \neq \text{Interrupt} \longrightarrow \text{ct-active } s) \cap is-subject <math>pas \circ \text{cur-thread}$ } call-kernel ev{ λ -. pas-refined pas}

We discuss the proof of these two theorems in the next section.

Via the refinement proof shown in previous work [12], both of these properties transfer to the C code level of the kernel. The guarantees that integrity pas st makes about user memory transfer directly, but the guarantees pas-refined pas and integrity pas st make about kernel-private state are mapped to their image across the refinement relation, which means we may lose some precision. The protection state of the kernel maps across the refinement relation precisely, the only difference between the model-level and C-level capability types being encoding.

The remainder of the system state does not translate so simply, but we contend that this does not matter. We envision the integrity theorem being useful mainly as a framing rule, with a component programmer appealing to the integrity theorem to exclude interference from other components and to the kernel model to reason about the component's own actions. In this case the programmer is interested not in the precise C state of the private kernel data, but about the related kernel model state. The integrity theorem will then provide exactly what is needed.

For proving confidentiality in the future, we may have to be more careful, because abstraction may hide sources of information that exist in the C system.

4 Proof and Application

4.1 Proof

The bulk of the proof effort was showing two Hoare triples for each kernel function: one to prove **pas-refined** as a postcondition, and one to prove **integrity**. These lemmas are convenient to use within the Hoare framework as we can phrase them in a predicate preservation style. In the case of the integrity predicate, we use a form of the Hoare triple (the left hand side of the following equality) which encapsulates transitivity and is easy to compose sequentially. This form is equivalent to the more explicit form (the right hand side) as a consequence of reflexivity and transitivity:

$$\forall P. (\forall st. \{ \lambda s. \text{ integrity } pas st s \land P s \} f \{ \lambda rv. \text{ integrity } pas st \}) = (\forall st. \{ \lambda s. s = st \land P s \} f \{ \lambda rv. \text{ integrity } pas st \})$$

The proof was accomplished by working through the kernel's call graph from bottom to top, establishing appropriate preconditions for confinement and integrity for each function. Some appeal was made to previously proven invariants.

The proof effort generally proceeded smoothly because of the strength of the abstraction we are making. We allow the subject to change arbitrarily anything with its label, we map most kinds of access to the **Control** authority, and we require anything to which the subject has **Control** authority to share the subject's label. These broad brushstroke justifications are easy to apply, and were valid for many code paths of the kernel.

For example, CSpace updates always occur within the subject. As preconditions for various functions such as cap-move and cap-insert we assert that the address of the CSpace node being updated has the subject's label and, for pas-refined preservation, that all authority contained in the new capabilities being added is possessed by the subject in the policy. These preconditions are properties about the static policy, not the dynamic current state, which makes them easy to propagate through the proof.

The concept of a static policy gave us further advantages. By comparison, we had previously attempted two different variations on a proof that seL4 directly refines the take-grant security model [15]. These proof attempts were mired in difficulties, because seL4 execution steps are larger than take-grant steps. In between the take-grant steps of a single seL4 kernel call, their preconditions may be violated because capabilities may have moved or disappeared, and so the steps could not be composed easily. This seemed particularly unfortunate considering that the take-grant authority model has a known static supremum. Comparing against this static graph instead yields something like our current approach.

Another advantage we have in this proof effort is the existing abstraction from the C code up to our kernel model. The cdt (capability derivation tree) and endpoint queues described already must be implemented in C through pointer linked datastructures. In C the cdt is encoded in prefix order as a linked list. When a subject manipulates its own capabilities it may cause updates to pointers in list-adjacent capabilities it does not control. In the abstract kernel model the cdt is represented as a partial map from child to parent nodes, making all of our subject's CSpace operations local. It would be possible to phrase an integrity predicate on the C level which allowed appropriate pointer updates, but we think it would be extremely difficult to work with.

The combined proof scripts for the two access control properties, including the definitions of all formalisms and the SAC example, comprise 10500 lines of Isabelle/HOL source. The proof was completed over a 4 month period and consumed about 10 person months of effort. Relative to other proofs about the kernel model this was rapid progress. Modifications to this proof have also been fast, with the addition of the cdt-integrity aspect of the integrity property being finished in a single day.

During the proof we did not find any behaviour in seL4 that would break the integrity property nor did we need to change the specification or the code of seL4. We did encounter and clarify known unwanted API complexity which is expressed in our policy wellformedness assumption. One such known problem is that the API optimisation for invoking the equivalent of a remote procedure call in a server thread confers so much authority to the client that they have to reside in the same policy label. This means the optimisation cannot be used between trust boundaries. An alternative design was already scheduled, but is not implemented yet.

We have found a small number of access control violations in seL4's specification, but we found them during a previous proof before this work began. These problems related to capability rights that made little sense, such as read-only Thread capabilities, and had not been well examined. The problem was solved by purging these rights.

4.2 Application

The application scenario of the integrity theorem is a system comprising trusted as well as untrusted components such as in the SAC example. Such scenarios are problematic for purely mandatory access control systems such as subsystems in a take-grant setting [15,7], because the trusted component typically needs to be given too much authority for access control alone to enforce the system's security goals (hence the need for trust). Our formulation of integrity enforcement per subject provides more flexibility. Consider a sample trace of kernel executions on behalf of various components in the SAC.

$$s_0 - - T - s_1 - - CTR - s_2 - RM - s_3 - R - s_4 - T - - T$$

The example policy in Fig. 3 satisfies our wellformedness condition for the components R, T, and CTR. It does not satisfy wellformedness for RM. This means, given the initial state s_0 , we can predict bounds for authority and state mutation using our theorems up to s_2 before RM runs. Since RM is the trusted component, we do not apply the integrity theorem, but reason about its (known) behaviour instead and get a precise characterisation of s_3 . From here on, the integrity theorem applies again, and so on. Suitably constructed, the bounds will be low enough to enforce a system wide state invariant over all possible traces which in turn, if chosen appropriately, will imply the security goal.

As hinted at in Sect. 2, at state s_2 , the RM component could use its excessive authority to reconfigure the system such that a new R is now connected to networks A and D. This setup would be incompatible with the policy applied to the transitions before, but a new policy reflecting the reconfiguration can be constructed that applies from there on. If that new policy and the reconfiguration are compatible with the system invariant we can conclude the security goal for a system that dynamically changes its high-level policy even if the integrity theorem itself assumes a static policy per kernel event.

For reasoning about such systems, it is convenient to lift wellformedness and therefore **pas-refined** to sets of allowed subjects instead of a single current actor. This means the same instantiation of the theorem can be applied and chained over traces without further proof. We have formulated and proved the lifted version, but omit the details here. They add no further insight.

The set versions of integrity and pas-refined are also useful in a more restricted, but common scenario where the kernel is employed as a pure separation kernel or hypervisor. In this scenario, all components would be considered untrusted, and there can be one system-wide policy that is wellformed for the set of all components. Wellformedness can be shown once per system and the subjective integrity theorem collapses to a traditional access control formulation of integrity that applies to all subjects.

4.3 Limitations

The limitations of the theorem as presented mostly reflect API complexities that we circumvented by making assumptions on the policy.

The strongest such assumption is that we require two components with a Grant connection to map to the same label. This is no more than what a traditional take-grant analysis amounts to [7], but in our subjective setting there would be no strong reason to forbid Grant from trusted to untrusted components if the authority transmitted is within policy bounds. The difficulties with this and with interrupt delivery were discussed in Sect. 3.2. Trusted components can still delegate capabilities via shared CNodes, which appear in our graph as Control edges. This is the approach taken by the RM component of the SAC.

Another limitation is that the theorem provides relatively coarse bounds. This is required to achieve the level of abstraction we are seeking, but it is imaginable that the bounds could be too high for a particular frame condition that is required. In this case, one could always fall back to reasoning about the precise kernel behaviour for the event under consideration, but it was of course the purpose of the theorem to be able to avoid this.

5 Related Work

Security properties have been the goal for OS verification from the beginning: the first projects in this space UCLA Secure Unix [19] and PSOS [8] (Provably Secure OS) were already aiming at such proofs. For a general overview on OS verification, we refer to Klein [11] and concentrate below on security in particular.

A range of security properties have been proven of OS kernels and their access control systems in the past, such as Guttman et al's [9] work on information flow, or Krohn et al [13] on non-interference. These proofs work on higher-level kernel abstractions. Closest to our level of fidelity comes Richards [17] in his description of the security analysis of the INTEGRITY-178B kernel in ACL2. Even this model is still connected to source code manually.

As mentioned, seL4 implements a variant of the take-grant capability system [15], a key property of which is the transitive, reflexive, and symmetric closure over all Grant connections. This closure provides an authority bound and is invariant over system execution. Similar properties hold for a broader class, such as general object-capability systems [16].

We have previously proved in Isabelle that the take-grant bound holds for an abstraction of the seL4 API [7,4]. The EROS kernel supports a similar model with similar proof [18]. However, these abstractions were simpler than the one presented here and not formally connected to code.

Compared to pure take-grant, our subjective formulation of integrity is less pessimistic: it allows certain trusted components, which are separately verified, to possess sufficient **Control** rights to propagate authority beyond that allowed by the policy.

6 Conclusions

In this paper, we have presented the first formal proof of integrity enforcement and authority confinement for a general-purpose OS kernel implementation. These properties together provide a powerful framing condition for the verification of whole systems constructed on top of seL4, which we have argued with reference to a case study on a real example system.

We have shown that the real-life complexity in reasoning about a generalpurpose kernel implementation can be managed using abstraction. In particular, our formalisation avoids direct reasoning about the protection state, which can change over time, by representing it via a separate policy abstraction that is constant across system calls. Integrity asserts that state mutations must be permitted by this policy, while authority confinement asserts that the protection state cannot evolve to contradict the policy.

This work clearly demonstrates that proving high-level security properties of real kernel implementations, and the systems they host, is now a reality. We should demand nothing less for security-critical applications in the future.

Acknowledgements

We thank Magnus Myreen for commenting on a draft of this paper.

This material is in part based on research sponsored by the Air Force Research Laboratory, under agreement number FA2386-10-1-4105. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright notation thereon. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of the Air Force Research Laboratory or the U.S. Government.

References

- J. Andronick, T. Bourke, P. Derrin, K. Elphinstone, D. Greenaway, G. Klein, R. Kolanski, T. Sewell, and S. Winwood. Abstract formal specification of the seL4/ARMv6 API. http://ertos.nicta.com.au/software/seL4/, Jan 2011.
- J. Andronick, D. Greenaway, and K. Elphinstone. Towards proving security in the presence of large untrusted components. In G. Klein, R. Huuck, and B. Schlich, editors, 5th SSV, Vancouver, Canada, Oct 2010. USENIX.
- D. Bell and L. LaPadula. Secure computer system: Unified exposition and Multics interpretation. Technical Report MTR-2997, MITRE Corp., Mar 1976.
- A. Boyton. A verified shared capability model. In G. Klein, R. Huuck, and B. Schlich, editors, 4th SSV, volume 254 of ENTCS, pages 25–44. Elsevier, 2009.
- D. Cock, G. Klein, and T. Sewell. Secure microkernels, state monads and scalable refinement. In O. A. Mohamed, C. Muñoz, and S. Tahar, editors, 21st TPHOLs, volume 5170 of LNCS, pages 167–182, Montreal, Canada, Aug 2008. Springer-Verlag.
- J. B. Dennis and E. C. Van Horn. Programming semantics for multiprogrammed computations. CACM, 9:143–155, 1966.
- D. Elkaduwe, G. Klein, and K. Elphinstone. Verified protection model of the seL4 microkernel. In J. Woodcock and N. Shankar, editors, *VSTTE 2008*, volume 5295 of *LNCS*, pages 99–114, Toronto, Canada, Oct 2008. Springer-Verlag.
- R. J. Feiertag and P. G. Neumann. The foundations of a provably secure operating system (PSOS). In AFIPS Conf. Proc., 1979 National Comp. Conf., pages 329–334, New York, NY, USA, Jun 1979.
- J. Guttman, A. Herzog, J. Ramsdell, and C. Skorupka. Verifying information flow goals in security-enhanced linux. J. Comp. Security, 13:115–134, Jan 2005.
- 10. T. Jaeger. Operating System Security. Morgan & Claypool Publishers, 2008.
- 11. G. Klein. Operating system verification—an overview. Sādhanā, 34(1):27–69, 2009.
- G. Klein, K. Elphinstone, G. Heiser, J. Andronick, D. Cock, P. Derrin, D. Elkaduwe, K. Engelhardt, R. Kolanski, M. Norrish, T. Sewell, H. Tuch, and S. Winwood. seL4: Formal verification of an OS kernel. In 22nd SOSP, pages 207–220, Big Sky, MT, USA, Oct 2009. ACM.
- M. Krohn and E. Tromer. Noninterference for a practical DIFC-based operating system. In *IEEE Symp. Security & Privacy*, pages 61–76, 2009.
- B. W. Lampson. Protection. In 5th Princeton Symposium on Information Sciences and Systems, pages 437–443, Princeton University, Mar 1971. Reprinted in Operat. Syst. Rev., 8(1), Jan 1974, pp 18–24.
- R. J. Lipton and L. Snyder. A linear time algorithm for deciding subject security. J. ACM, 24(3):455–464, 1977.
- T. Murray and G. Lowe. Analysing the information flow properties of objectcapability patterns. In 6th FAST, volume 5983 of LNCS, pages 81–95, 2010.
- R. J. Richards. Modeling and security analysis of a commercial real-time operating system kernel. In D. S. Hardin, editor, *Design and Verification of Microprocessor* Systems for High-Assurance Applications, pages 301–322. Springer-Verlag, 2010.
- J. S. Shapiro and S. Weber. Verifying the EROS confinement mechanism. In *IEEE Symp. Security & Privacy*, pages 166–181, Washington, DC, USA, May 2000.
- B. J. Walker, R. A. Kemmerer, and G. J. Popek. Specification and verification of the UCLA Unix security kernel. *CACM*, 23(2):118–131, 1980.



University of New South Wales School of Computer Science and Engineering

Formal Functional Specification of a Security Critical Component in Isabelle/HOL

Author: Nelson Billing

Supervisor: June Andronick Assessor: Gerwin Klein

October 18, 2011

A Thesis submitted in partial fulfillment of the requirements for the degree of Bachelor of Engineering (Software Engineering)

Acknowledgements

I would like to thank the following persons: Ning Yu, June Andronick, Gerwin Klein, Thomas Sewell, David Greenaway, Bernard Blackham, Rafal Kolanski, David Cock, and Matthias Daum. This list is not intended to be exhaustive, and is in no particular order.

Contents

1	Intro	Introduction 6			
2	Bacl	kground 9			
	2.1	Forma	l verification	9	
	2.2	The se	L4 Microkernel	10	
	2.3	seL4's	Formal Proof of Correctness	11	
	2.4	The Se	ecure Access Controller	12	
	2.5	The R	outer Manager	14	
	2.6	Access	s control in seL4	15	
3	Des	ign		16	
	3.1	Route	r manager design	17	
			4 design	18	
		3.2.1	ARM & x86	19	
		3.2.2	Generic libseL4	19	
	3.3	System process design			
		3.3.1	System process from seL4 refinement	19	
		3.3.2	Modifications	20	
		3.3.3	Generic system process	21	
4	Imp	nplementation 2			
	4.1	4.1 Router Manager Functional specification		22	
		4.1.1	The sac_abstract_state	23	
		4.1.2	Router Manager program instructions	24	
		4.1.3	Router Manager program structure	26	
		4.1.4	Example: <i>Teardown Router</i>	27	
		4.1.5	Example: Setup Router	28	

		4.1.6 In-depth <i>instr</i> examples	29			
		4.1.7 Internal helper function examples	32			
	4.2	2 libseL4 specification				
		4.2.1 The <i>tc_monad</i>	34			
		4.2.2 The <i>thread_abstract_context</i>	34			
		4.2.3 libseL4 specification structure	34			
		4.2.4 In-depth libseL4 examples	35			
		4.2.5 libseL4 helper function examples	37			
		4.2.6 Lifting from <i>tc_monad</i> to <i>rm_monad</i>	38			
	4.3	System process	39			
		4.3.1 The global_automaton	39			
		4.3.2 The <i>ADT_I</i>	41			
		4.3.3 Kernel specification	42			
		4.3.4 User specification	42			
_	6					
5			44			
	5.1	Contributions				
	5.2	Future Work	44			
Α	A Router Manager Functional Specification					
	A.1	The Router Manager Functional specification program	46			
	A.2	The sac_abstract_state	47			
	A.3	The rm_monad and Router Manager execution	48			
	A.4	Router Manager helper functions	49			
	A.5	Router Manager instructions	50			
	A.6	Router Manager program	57			
-	1		~			
В		1	61			
	B.1	libseL4 Interface				
	B.2		65			
	B.3		66			
	В.4	libseL4 Types	66			
С	Constant Bracass					
C	Svst		73			
	-	tem Process	73 73			
	-	tem Process Router Manager Functional specification system process	73			
	C.1	tem Process Router Manager Functional specification system process C.1.1 Lifting from rm_monad to user_state transitions	73 73			
	C.1	tem Process Router Manager Functional specification system process C.1.1 Lifting from rm_monad to user_state transitions User state translations	73			

C.4	The system process	75
C.5	System process	75
C.6	The global_automaton	76

Chapter 1

Introduction

Computers are extremely pervasive today, and so of course is software. Software is notoriously difficult to get right, for the most part developers just accept uncertainty in a piece of software to some degree.

In software applications with very little fault tolerance, the assurance given by standard software testing is unsatisfactory and so there are a number of approaches to achieving higher levels of assurance. One obvious approach is just to do more testing, the more testing you do the closer to complete coverage of inputs and therefore less likely that bugs will slip through.

Alternatively there are formal verification and static code analysis techniques which can deliver stronger claims about correctness but which are typically more costly than testing. Therefore formal methods techniques are applied where the need for assurance justifies the cost.

The goal of the trustworthy embedded systems project at NICTA, where this thesis project is taking place, is to build systems we can really trust, by providing formal proofs about their security, safety and reliability. The first step in this vision was the creation of the seL4 microkernel, a general purpose microkernel of the L4 family, along with a formal proof of its functional correctness[3]. The proof of functional correctness was completed in 2009 and since then work has been going on to build larger systems on top of the seL4 microkernel, and prove specific, targeted properties about those systems. This thesis project is a part of these next steps of the vision. The idea being how can we use this verified kernel to prove properties about user programs?

The Secure Access Controller (SAC) is a case study in verification of systems built on seL4. The role of the Secure Access Controller is a secure router between two networks which must be kept isolated- meaning that the targeted property we want to hold about the SAC is absence of data flow between the networks.

The trustworthy embedded systems project already have a very high level security model of the

SAC and a proof of absence of data flow on that model[1]. The aim of the larger project, which this thesis is a part of, is connecting the SAC Security model to the implementation of the SAC so that the proofs on the Security model hold for the source code as well. Connecting the two up means showing that the SAC source code has the same behaviour as the security model used for the proof.

We decided to show refinement via an intermediate step. The SAC's architecture consists of only one trusted component, called the Router Manager and, as we will show, the behaviour of *untrusted* components can safely be underspecified as arbitrary for the purposes of verification. Therefore the intermediate step would be a Functional specification of the Router Manager component of the SAC.

The Router Manager component must make calls into the seL4 kernel to perform its function, such as to create and setup new threads, and to receive control messages. In the C implementation of the Router Manager component, calls to the kernel are made by calling functions from a C library built for seL4, called libseL4. Instead of effectively 'inlining' the libseL4 code into the Router Manager Functional specification, we would like to implement a libseL4 specification. It is desirable for us to implement the libseL4 specification such that it can be used in other seL4 user program specifications. For this reason, we stress that the libseL4 specification must be *generic*.

With the Functional specification done, we would need a framework for proving the new Functional specification refines the Security model. The framework would preferably be general in the specification of user behaviour so that it can be reused in showing refinement between other sel4 user program specifications. In other words, we would like to be able to plug any user program specification into the system process framework and reason about the resulting system process.

Writing the Router Manager Functional specification was the main objective of this thesis, with the aim of including a generic libseL4 specification. A stretch goal was to start on a framework for the refinement proof.

We completed the Functional specification, which involved specifying a subset of the libseL4 library alongside it. The libseL4 subset was implemented in a generic fashion so that it can be used in specifying other user programs (and we would expect that most or all other seL4 user programs do use libseL4 in their C implementation). We also implemented the framework for stating refinement between seL4 user programs, along with an example of its use with the Functional specification.

This document firstly presents the necessary background in formal verification, seL4, and the Secure Access Controller in Chapter 2. Following that is Chapter 3 which presents the design of the Router Manager Functional specification, libseL4 specification, and refinement framework, with the aim of explaining the choices in implementation and providing some much more detailed background.

Chapter 4 is structured around the three main parts of the project that were implemented. In each section, details from the Isabelle/HOL theories are presented and discussed in order to show the reader the 'hows' of the implementation. The final chapter of this document, Chapter 5, summarises the contributions made in this project and offers the experience gained during the work, as well as what remains to be done in future work.

Chapter 2

Background

This chapter is intended to cover some general background on topics related to this thesis project which may be helpful to readers. Section 2.1 gives an overview of formal methods techniques, the particular technique used in this thesis, and a justification of our choice. Next, Section 2.2 will introduce the seL4 microkernel and give some background for readers who are unfamiliar. Section 2.3 introduces the techniques used in the seL4 formal verification project which guided the approach in this project. After that, Section 2.4 presents some background on the target of the specification project, the Secure Access Controller, its components and its design. As a followup, Section 2.5 briefly discusses the central component of the SAC, the Router Manager. The last section of this chapter, Section 2.6, gives some background information on how access control is handled in seL4 and how this guides the componentization of the SAC.

2.1 Formal verification

There are several broad approaches to formal verification with different levels of expressivenessin terms of the properties about programs- and automation.

Model checking, where we create a model of our system and then use an automated tool to explore its state space, gives us a lot of expressiveness in talking about properties of a program. While verification is fully automatic, it is not feasible for all programs because of the so-called *state explosion* problem. Besides, the technique does not involve showing refinement from source code to the model, so it does not prove anything directly about the source code.

Static code analysis techniques can be used to automatically detect patterns of programming mistakes- like memory leaks or out-of-bounds memory access- so that specific strong claims can be made about programs. The nature of automatically checking source code for 'mistakes' does

not lend itself to proving arbitrary program-specific properties.

Theorem proving enables us to reason formally about a logical model of the program that can be strongly linked to the program itself using a formal semantics of the programming language of the program. This technique allows us to phrase arbitrary properties about a program and prove them at the source code level. Unlike model checking, state space is not immediately important because we are not bound to prove by exhaustion. The price we pay with theorem proving is that each stage of the process requires large amounts of human work along with mechanical assistance.

The formal verification technique which was chosen for this thesis project is theorem proving. The reasons for this choice are the same reasons that it was chosen for the seL4 verification. First the large state space of seL4 and the SAC makes automatic state space exploration impractical. Second, our goal is a code level guarantee. And finally, we are aiming at specific properties that need to be phrased in an expressive language.

Verification of the seL4 kernel was carried out using the Isabelle theorem prover [6], and the models of the seL4 kernel are written as Isabelle/HOL theories. Since this thesis project aims at leveraging the proofs of the seL4 kernel, it is natural to continue using Isabelle/HOL, which also happens to be the theorem prover I have experience with.

2.2 The seL4 Microkernel

The seL4 microkernel is a L4 family microkernel, with capability-based security [2]. seL4 is the first general purpose OS kernel with formal proof of its functional correctness.

seL4 provides the usual L4 abstractions: threads, address spaces, and inter-process communication (IPC). Threads, address spaces, IPC and all other (more esoteric or architecture specific) abstractions are represented by *kernel objects*. Authority over *kernel objects* is conferred via capabilities. In fact, all kernel calls in seL4 are invocations on a capability- for instance mapping a frame into an address space is an invocation on a capability to the frame.

Capabilities in seL4 serve a similar function to Linux file handles. However they are meant to implement much stricter access control policy. Capabilities are kept *for* threads, in their respective *CSpaces* which are conceptually similar to address spaces. The capabilities themselves are unforgeable and threads are only ever given pointers to the *CSpace* and then into *slots* inside it, so as to prevent any tampering with the capability kernel object. There are strict controls on granting capabilities over IPC, and threads may only copy or move capabilities between or inside *CSpaces* indirectly via the kernel.

A CSpace in the kernel is made up of a hierarchy of CNodes. A CNode is a kernel object which

serves as a container for capabilities. The hierarchy is built by placing capabilities to one *CNode* inside another *CNode*, thereby allowing addresses to be 'walked' from a root *CNode*, through subsequent *CNodes* to some particular capability in that *CSpace*.

Rights to create new kernel objects in system memory are represented by *Untyped* capabilities, each having a start address and size to indicate which region of memory it covers. When a user wants to create a new kernel object, they must *retype* some part of an *Untyped* capability they hold to do so.

As kernel objects or new capabilities are created using the authority of a capability, the kernel keeps track of which objects and capabilities are *children* of that capability. The counterpart to *retyping* an *Untyped* capability is to call *revoke* on that capability. When *revoke* is performed on the original *Untyped* capability, all of the *child* capabilities are also revoked, and then destroyed. Any kernel objects associated with these capabilities are destroyed in the process.

Of course, when we speak of the user performing some action at the kernel objects or capabilities, or performing inter-process communication, what we mean is that the user is asking the kernel to perform this action on their behalf. Readers will be familiar with Unix syscalls and filehandles, and seL4 syscalls and capabilities roughly correspond to these respectively.

2.3 seL4's Formal Proof of Correctness

The seL4 microkernel makes a part of the trusted computing base for the SAC, and in that sense this project is building directly on the formal verification of functional correctness of seL4. Functional correctness here meaning that the source code of seL4 implements the desired specification. The process of verifying seL4 is well documented and [5][3][7] are instructive to the techniques used.

The first step in proving functional correctness is to write the formal specification of the program. This describes *what* the kernel is supposed to do, whereas the implementation describes *how* it does it. In this case the formal specification of the seL4 kernel is an Isabelle/HOL theory termed the *abstract specification*. The abstract specification has an abstract view of the seL4 kernel state and is nondeterministic. It is written in the functional style of Isabelle/HOL programs and uses nondeterministic and deterministic state monads.

The next step is to prove that the source code of the program actually implements the specification. The way this was done for seL4 is the *implementation*, which is a translation of the seL4 source code into Isabelle/HOL, is shown to be a refinement of the *abstract specification*.

Since the refinement relation is transitive the refinement from the *implementation* to the *abstract specification* can be completed in two steps: from the *implementation* to an *executable specification*

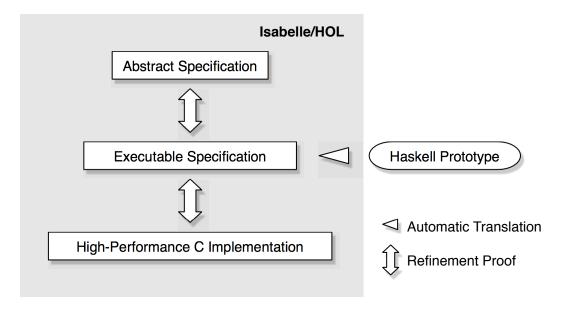


Figure 2.1: Refinement steps in seL4 refinement [3]

and from the executable specification to the abstract specification, as illustrated in Figure 2.1.

A concise argument for the soundness of the techniques used in the refinement is given in [5]. Each of the two refinement steps is proved by a forward simulation between two of the processes, the general shape of these forward simulations is illustrated in Figure 2.2. The idea of forward simulation is to show that given a *state relation* between two processes, when the processes are started in some related state and react to the same event they will end in some related states. Forward simulation is proven by showing an equivalent property termed *correspondence*, which is defined between terms in the specifications rather than *events*.

2.4 The Secure Access Controller

Ultimately this project is a case study in proving correctness of a software system. The system used is a Secure Access Controller (SAC), a router between a trusted user terminal and a number of untrusted networks, as illustrated in Figure 2.3. For simplicity and without loss of generality we will only consider two such untrusted networks.

The best way to explain the role of the SAC is to imagine a user who needs access to two different computer networks. The user is not allowed to connect to both networks at once. The administrators of these two networks trust the user not to save information to pass it on, but do not trust his network card/TCPIP stack/OS to keep their network isolated from other networks. The role

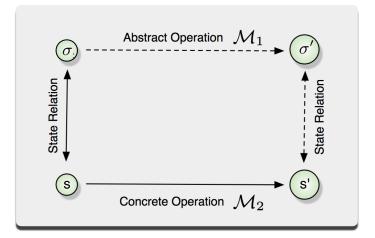


Figure 2.2: Forward simulation [3]

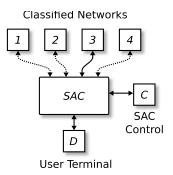


Figure 2.3: The SAC routes between a user's terminal and 1 of *n* classified networks [1]

of the SAC here is to sit between the user terminal and the two outside networks and connect the user to one network at a time. To satisfy the administrators we need to show that the SAC won't leak data between the networks.

The SAC needs a component responsible for network routing, including network card drivers. This component will be called the Router. It also needs a web interface where the user can request a network switch, this will be called the SAC Controller. These two components represent a significant amount of code, much larger than seL4. Proving non-trivial properties about that much source code is infeasible, so the guiding design principle of the SAC was a minimal trusted computing base (TCB) in order to simplify the proof process.

Our system of course needs to trust the behaviour of the seL4 microkernel- for which we have a proof- but the SAC program can also be split into several *components*, some being in the trusted

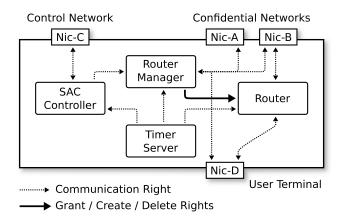


Figure 2.4: High level component breakdown of the SAC design [1]

computing base and others not. The logical way to split the program into components is to consider the least privilege each part of the program needs in order to function, so that components individually have a small subset of the capabilities that the SAC has as a whole.

The targeted security property we focus on in this thesis project is absence of data flow between networks, so the idea is to pinpoint components which can access both networks. For example, the Router is given access to network interfaces to the user terminal and *only one* of the outside networks. It was implemented by a Linux instance comprising millions of LOC. Therefore we place it outside of our TCB. To ensure we do not need to trust the Router we destroy it whenever we switch from one outside network connection to another. This critical operation is done by a trusted component, called the Router Manager.

The Router Manager has complete access to both external networks and thus is able to break our security property. This places the Router Manager inside our TCB. The remaining components are the SAC Controller, which is given access only to the control network, and a timer server. As we can see in Figure 2.4, only the Router and Router Manager components have access to the outside networks. Therefore given some reasonable assumptions about seL4 enforcing access control, discussed in Section 2.6, the only trusted component we need to verify is the Router Manager.

2.5 The Router Manager

The Router Manager (RM) component of the SAC system is the core logic of the system and the one *trusted* component. The Router Manager sits in an infinite loop and waits for instructions from a Control Network interface. Once it receives a message it can either just teardown the

Router component of the SAC or teardown the Router component and then setup a new Router to either Network A or Network B. Once the Router Manager has finished responding to the control message it will return to waiting for another message. This process is discussed in more detail in Subsection 4.1.3.

Setting up the Router component involves creating a Router *CSpace*, a Router virtual address space, a Router Inter-Process Communication buffer, a Router Thread Control Block, and then connecting them all together and activating the new thread. Regions of memory such as the Router's program text and some free space need to be mapped into the Router's virtual address space. Additionally, the capabilities to either the Network A interface or Network B interface, and some *Untyped* capabilities, need to be transfered into the Router's *CSpace*. This process is explained in more depth in Subsection 4.1.5.

Tearing down the Router component involves calling *revoke* (see Section 2.2) on all *Untyped* capabilities which were used in setting up the Router, as well as 'sanitising' whichever network interfaces the Router was connected to (Netword D and either Network A or Network B). This process is explained in more depth in Subsection 4.1.4.

2.6 Access control in seL4

In order to know that the untrusted SAC components cannot interfere with our security property we need to know that seL4 is enforcing an access control policy on these processes. For example we need to know that the seL4 kernel will not allow the Router to access both outside networks with just a capability to one, or move capability pointers around in the Router Manager's CSpace. Showing that seL4 enforces this access control policy is ongoing in the verification team at NICTA. One half is shown in the recently completed Integrity theorem [4], while the remaining half is the work in progress Confidentiality theorem. Roughly, the Integrity theorem shows that processes are bounded in writes and Confidentiality is intended to show that processes are bounded in reads.

Chapter 3

Design

The aim of this chapter is to step through the design of the Router Manager Functional specification, as well as the system process that we want to be a combination of the kernel and Router Manager. The Router Manager Functional specification is meant to specify the behaviour of the Router Manager component of the Secure Access Controller, that is, it gives the set of all seL4 user state transitions that the Router Manager can make.

The Router Manager needs to perform system calls in setting up or tearing down other processes or receiving messages. In the C implementation of the Router Manager the system calls are done via the libseL4 library and so it is necessary to specify some subset of libseL4 in order to have a full specification of the Router Manager that can interact with a specification of the seL4 kernel.

With libseL4 giving the interface to the kernel, we also need a kernel specification, so that the system process can 'trap' into the kernel after a syscall is set up by the user process. Specifications of the seL4 kernel at varying levels of abstraction already exist.

The final piece of the puzzle in defining the system process is a framework that will combine the transitions of the user specification (the Router Manager Functional specification) and the seL4 kernel. This requires definitions to translate or 'lift' between views of system state and a definition of the overall system transitions given the user and kernel transitions. The latter is referred to in this document as an *automaton*.

This chapter first discusses the design of the Router Manager in Section 3.1, then the design of libseL4 in Section 3.2, and finally Section 3.3 presents the design of the system process.

3.1 Router manager design

The purpose of the Router Manager Functional specification, is to provide a less jarring refinement step for both the security model and implementation. Thus the suitable level of abstraction for operations and data in the Router Manager Functional specification would lie somewhere between the implementation and security model. To contrast with the verification of seL4: the high level security model of the SAC is far more abstract than the *abstract specification* of the seL4 kernel, mentioned in Figure 2.1. seL4's *abstract specification* is a full functional specification of the kernel, whereas the SAC's security model abstracts away everything that is not needed for the security analysis. The Router Manager Functional specification is intended to employ a similar level of abstraction to the seL4 *abstract specification* and describe the RM's functional behaviour, as shown in Figure 3.1.

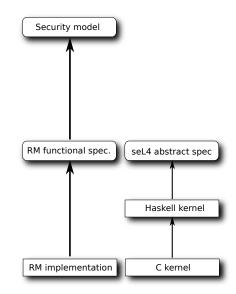


Figure 3.1: Side-by-side comparison of SAC and seL4 refinement artefacts.

The Router Manager Functional specification of course will need a program state capturing the state we want to project out of the security model and implementation's states. The content of the state contains most of the internal state of the RM implementation in a more abstract form.

There is a fundamental mismatch between the security model's RM program and the RM implemention in that the first is unstructured code, branching with jumps, while the second has the program structure of a C program- *if*, *for*, *while*, etc. A middle ground between structured and unstructured code is unhelpful, so the Router Manager Functional specification represents programs as unstructured code.

datatype sysOPs =				
	SysRead	cap		
I	SysWrite	cap	bool	
I	SysCreate	cap		
I	SysGrant	cap	cap	
I	SysRemoveSet	cap	(cap	set)
1	SysDelete	cap		

Figure 3.2: Definition for kernel calls in Router Manager security model.

RM steps can be divided into internal processing and calls into the seL4 kernel. In the security model kernel calls are represented by an abstract instruction datatype, shown in Figure 3.2, and the only internal processing is the jump instructions which move control to some other position in the program non-deterministically.

The RM implementation takes mostly internal steps punctuated by kernel calls to set up or communicate with other processes. Details on how the Router Manager Functional specification has been defined are given in Chapter 4. The kernel calls involve setting up registers and trapping into the kernel; this process is wrapped in libseL4 library functions which the RM's implementation uses, and which are discussed in the next section.

3.2 libseL4 design

As mentioned in Section 3.1, the router manager program uses the libseL4 library as an interface with the seL4 microkernel. Therefore, to fully specify the router manager program we need a specification of the libseL4 library. The router manager needs to cross the kernel boundary to create a new thread, transfer capabilities, retype untyped memory into new kernel objects, or to map memory into another address space et cetera, so all of these tasks will involve executing code from libseL4. Unfortunately no such specification existed, requiring me to specify the parts of libseL4 that are used by the router manager.

The libseL4 library provides a simple interface to the kernel, as it is described in the seL4 manual, in terms of a set of *'invocations'* for each capability type. Internally the kernel decodes calls made to it and then becomes aware of what the *invocation* it is handling, but the user side of the interface does not have functions for each *invocation*, and arguments to the kernel calls and magic numbers must be put into registers for the kernel to decode. Thus libseL4 is useful as providing an abstraction of the real calling mechanism, as well as marshalling and copying data that the user passes as arguments.

The implementation of the libseL4 specification is presented in detail in Section 4.2.

3.2.1 ARM & x86

The architecture specific code in the router manager comprises just calls to the kernel interface functions, so selecting the architecture that the router manager specification would assume is essentially choosing the version of libseL4 to use. Since the existing router manager program was written for the x86 architecture whereas an ARM- but not x86- specification of the seL4 kernel exists, I had to make a choice between the two architectures for my libseL4 specification.

Inevitably I chose to specify the ARM libseL4 implementation, since the seL4 kernel specification is overwhelmingly larger and more complex than the router manager program. Before the refinement from the Router Manager Functional specification to the C code, the router manager C code will need to be ported to ARM. Given the small amount of architecture specific code in the router manager this should not alter the overall structure of the program.

3.2.2 Generic libseL4

Most of the motivation for this project lies in completing a case study in verifying software, specifically what sort of leverage we can net by designing that software as a seL4 user program. With that in mind, developing the specification for a subset of the libseL4 library needed to be done in such a way that it would be agnostic to a user program specification importing it. This way any other seL4 user program could be specified in the fashion most appropriate to that program, and where it would interface with the kernel the libseL4 specification could describe that behaviour.

Initially I identified two key features that a *generic* libseL4 specification would need. Firstly the state that libseL4 functions operate over would need to be a subset of any seL4 user program's state. Secondly the libseL4 specification should introduce all the types that the specification of a seL4 user program might use when communicating with seL4.

3.3 System process design

Given the Router Manager security model and the functional specification we have most of the components we need to begin showing that the C code for the router manager refines the security model, but still remaining is the framework for showing either of the refinement steps.

3.3.1 System process from seL4 refinement

As discussed in Section 2.3 the method of program verification for the seL4 kernel was refinement from an abstract specification to source code, via an *executable* specification. In that case refinement was shown by defining a process, representing the kernel and user reacting to events, for each kernel specification. What is shown, then, is that the behaviour of one process refines the other.

To this end, a somewhat generic kernel-system automaton already exists. Roughly speaking; a non-deterministic kernel specification can be given which dictates *kernel mode* behaviour, and *user mode* behaviour is unspecified. The three *modes* under which the system operates are shown in Figure 3.3. Separate kernel and user state is recognised, where the user state will be a projection on the kernel state.

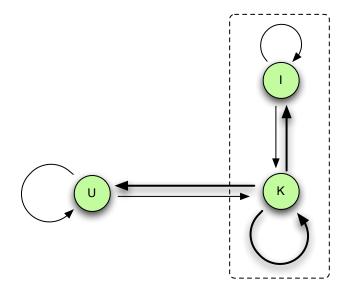


Figure 3.3: System transitions between *Idle*, *User*, and *Kernel* modes.

3.3.2 Modifications

In adapting this framework for the Router Manager refinement there were a number of immediate questions: which kernel specification would be used, should the user state be expanded, how should the kernel react to events, and should the automaton treat a kernel to user transition differently than a user to user transition?

The original intention in designing the RM Functional specification was to try for a level of abstraction similar to that in the seL4 Abstract specification, therefore it seemed logical to choose the seL4 Abstract specification as the kernel specification in the RM Functional specification system process.

User state needed to be expanded so that an abstract program counter called the instruction pointer

could be maintained for the active user process when the system is in *user mode*. The *instruction pointer* is necessary so that the user specification can control its own program flow and resume at the correct place when returning from a kernel call or being scheduled.

The system process updates the abstract *instruction pointer* on a *kernel mode* to *user mode* transition, using the program counter placed in a register by the kernel. A user specification is expected to update its *instruction pointer* in transitions, except when setting up a software interrupt since the kernel will update the program counter in a register when servicing the kernel call. This also answers the fourth system process design question above.

3.3.3 Generic system process

Being that this project is aiming to test-run seL4 user program verification, it is important that the newly modified system process is *generic* enough that it can represent any user system running on top of seL4. Specifically it should have a view of user state which will be a superset of any user program's observable state. The generic system process framework implementation is discussed in Section 4.3.

Chapter 4

Implementation

The purpose of this section is to describe in detail how it is that the Router Manager Functional specification was implemented. As discussed in Chapter 1, the original goal of this thesis project was functional specification of the Router Manager program. This goal is complete, and the details of its implementation are contained in Section 4.1.

Also discussed was that a stretch goal, or future work, from this project could be specifying the libseL4 library and/or creating a framework for proving refinement between the Router Manager Security model and the Functional specification. The parts of the libseL4 library which are used by the Router Manager program are specified, and in a generic manner which would allow them to be used, covered in Section 4.2. The foundation of a generic user refinement framework will be a user and kernel system automaton, I implemented this as well, covered in Section 4.3.

4.1 Router Manager Functional specification

The aim of the Router Manager Functional specification is to create a specification for the Router Manager which is more abstract than the Router Manager C implementation, but more concrete than the Secure Access Controller Security model. Structurally we attempted to mirror the C implementation and Security Model as often as possible.

This chapter will present the implementation of the Router Manager Functional specification. Subsection 4.1.1 discusses the Functional specification's state, Subsection 4.1.2 presents the way program transitions are represented in the Functional specification, Subsection 4.1.3 gives the structure of the Router Manager program, and the remaining subsections give examples from the Functional specification's Router Manager program.

The examples in the later subsections start from an instruction in the top level Functional spec-

ification Router Manager program and then expand more definitions in subsequent subsections until they reach the kernel barrier. Then the examples in Section 4.2 regress the same tasks further. The entire implementation of the Router Manager Functional specification is too big to discuss in detail in this report, but it is given in full in Appendix A.

4.1.1 The *sac_abstract_state*

record sac_abstract_state =
 ip :: seL4_Word
 untyped_items :: seL4_CPtr list list
 free_items :: seL4_Word list
 total_items :: seL4_Word list
 info :: seL4_MessageInfo
 message :: seL4_Word
 var_i :: nat
 RM_RWGC_NIC_A_cap :: seL4_CPtr
 RM_RWGC_NIC_B_cap :: seL4_CPtr
 RM_RWGC_NIC_D_cap :: seL4_CPtr
 client_vspace :: seL4_CPtr
 virtualirq_endpoint :: seL4_CPtr
 command_endpoint :: seL4_CPtr

Figure 4.1: Definition of the *sac_abstract_state*

The *sac_abstract_state* is an abstract representation of the Router Manager's user memory view. In terms of the Router Manager C program, this would correspond to allocated memory, local, and global variables. As mentioned in Section 3.3, the Router Manager program is expected to update its own *instruction pointer*, so this is contained inside the *sac_abstract_state*.

Figure 4.1 shows the *sac_abstract_state* global and local variables from the C implementation of the Router Manager. *ip* is the abstract *instruction pointer*, *info*, *message*, and *var_i* are local variables. The remaining ones are global variables.

untyped_items contains pointers to all of the *Untyped* capabilities the Router Manager has, *free_items* contains how many *Untyped* capabilities of each size are free, and *total_items* contains how many of each size *Untyped* the Router Manager has in total.

RM_RWGC_NIC_I_cap contains a pointer to the Router Manager's all-rights (**R**ead **W**rite **G**rant **C**reate) capability to Network I interface.

client_vspace, virtualirq_endpoint, and *command_endpoint* are pointers to capabilities which are set up by a booter initially and which we just pass around.

4.1.2 Router Manager program instructions

The Router Manager program is represented by a list of instructions. The instructions are themselves defined as functions across the Router Manager's abstract state and optionally producing an *event* indicating if a kernel event was generated by the instruction. The type of instructions is *(event option) rm_monad.*

The *rm_monad* type represents a computation which can modify the RM program state, shown in Figure 4.2. *nondet_monad* here is from some non-deterministic monad infrastructure written by Gerwin Klein and is defined in the usual way, with *bind*, *return* and sugar like *do*, *;* and \leftarrow . Functions which modify the RM program state can be defined independently of an *event* and then used in the definitions of different instructions. Functions which do not modify the RM program state- that would be so-called *pure* functions- can be defined as ordinary Isabelle/HOL functions over their input and output types.

As shown in Figure 4.2 the *rm_monad* state comprises a *sac_abstract_state* and *thread_abstract_state*. The former being the state which is considered specific to the Router Manager program implementation, and the latter being the state which is common to any seL4 user program. The reason that a clean separation between the two is desirable will be explained in Section 4.2. The *thread_abstract_context* will be explained in Subsection 4.2.2.

types	'a rm_monad	=	((thread_abstract_context \times sac_abstract_state), 'a) nondet_monad
types	instr	=	(event option) rm_monad
types	program	=	instr list

Figure 4.2: Definition of program terms in Router Manager functional specification.

rm_monad is a non-deterministic state monad. This is because, while the Router Manager Functional specification *is* deterministic, there is no reason for the generic seL4 user program verification machinery to assume determinism. In fact, the SAC Security model is non-deterministic, and so this choice also means we can use the same machinery to build a process for the SAC Security Model.

For computations of the type *unit rm_monad*, we wrap the non-event-generating computation in a function to increment the *instruction pointer* and return *None*, signifying no event.

As shown in Figure 4.3 the k_instr wrapper is the identity and is included to balance the u_instr wrapper, where k_instr is used on instructions which trap into the kernel and u_instr is used on instructions which trap into the kernel do not need to increment the

instruction pointer, as the kernel will do this, nor do they need to be wrapped in a *None,* since they are guaranteed to return *Some* (*SyscallEvent sy*).

```
definition k_instr ::
    instr ⇒ instr
    where
    k_instr = id

definition u_instr ::
    'a rm_monad ⇒ instr
    where
    u_instr m =
    do m;
        incr_pc;
        return None
    od
```

Figure 4.3: *u_instr* and *k_instr* definitions.

The definition of an *event* in Figure 4.4 comes from the seL4 *Executable specification* (see Figure 2.1) of the seL4 microkernel. This type captures all of the possible events that the seL4 kernel responds to. For the purposes of the Router Manager functional specification we need only consider *SyscallEvent*, since a user program will cannot intentionally spawn the other events. On ARM-our chosen architecture- this corresponds to a software interrupt.

```
datatype event =
    SyscallEvent syscall
    UnknownSyscall nat
    UserLevelFault machine_word machine_word
    Interrupt
    VMFaultEvent vmfault_type
```

Figure 4.4: Definition of *event*.

Again, the definition of *syscall* in Figure 4.5 comes from the seL4 *Executable specification*. The Router Manager functional specification makes libseL4 calls which only use *SysCall* and *SysWait*, however a full libseL4 specification would potentially use all of them. For clarity: every Router Manager functional specification instruction returns either *None*, *Some* (*SyscallEvent SysCall*) or *Some* (*SyscallEvent SysWait*).

```
datatype syscall =
    SysSend
    SysNBSend
    SysCall
    SysWait
    SysReply
    SysReplyWait
    SysYield
```

Figure 4.5: Definition of syscall.

As discussed in Section 3.1, the Router Manager program is an unstructured list of instructions. Branching and conditional execution in the Router Manager program is achieved by 'jumping', we define a *jmp* and *cnd_jmp* in Figure 4.6.

```
definition jmp ::
    \texttt{seL4\_Word} \ \Rightarrow \ \texttt{instr"}
   where
   jmp n =
       do modify (\lambda(tc, rm). (tc, ip_update (\lambdax. x + 1) rm));
           return None
       od"
definition cnd_jmp ::
   (\texttt{(thread_abstract\_context} \times \texttt{sac\_abstract\_state}) \Rightarrow \texttt{bool}) \Rightarrow \texttt{seL4\_Word}
\Rightarrow instr"
   where
   cnd_{jmp} P n =
        (\lambda s. (if (P s))
                then
                   jmp n
                else return None) s)
```

Figure 4.6: Definitions of *jmp* and *cnd_jmp*.

4.1.3 Router Manager program structure

At the top level there are two parts to the functional specification Router Manager program. Firstly it waits for a message from the Control Network. Secondly it executes one of three branches of a switch on the received message.

The three branches are 'teardown Router', 'create Network A Router', and 'create Network B Router'. The first branch revokes all kernel objects which were created in setting up the Router thread, thereby destroying the Router and erasing anything the Router wrote to memory. The second branch also revokes all Router kernel objects, and then creates a new Router thread and gives it access to Network A. The third branch does the same as the second branch, except it gives access to Network B instead of Network A.

Once a branch has finished, the Router Manager program returns to waiting for a message from the Control Network. Figure 4.7 is illustrative to the structure.

At the stage that this part of the program runs, the Router Manager has just resumed from waiting on a message from its Control Network interface.

As covered in Section 2.4 and Subsection 4.1.3, the Router Manager's behaviour is quite simple.

Figure 4.7 shows precisely how the Router Manager algorithm is represented in the Router Manager functional specification. Also shown is the corresponding C code, which is very similar and should aid in understanding the flow of the unstructured code.

```
C code
                                                     RM Line
                                                                 RM Instruction
info = seL4_Wait(command_endpoint, NULL);
                                                             0
                                                                 k_instr wait_for_message
message = seL4_GetMR(0);
                                                             1
                                                                 u_instr get_message
switch (message) {
                                                             2
                                                                 cnd_jmp (\lambda(tc,rm). message rm = 0) 4,
  case 0:
                                                             3
                                                                 jmp 13,
     (Teardown Router)
                                                                 (Teardown Router)
     break;
                                                            12
                                                                 jmp 86
                                                            13
                                                                 cnd_jmp (\lambda(tc, rm). message rm = 1) 15,
  case 1:
                                                            14
                                                                 jmp 50
     (Teardown Router)
                                                                  (Teardown Router)
     \langle \text{Setup Router (NIC A)} \rangle
                                                                 \langle \text{Setup Router (NIC A)} \rangle
                                                            49
     break;
                                                                 jmp 86
                                                            50
                                                                 cnd_jmp (\lambda(tc, rm). message rm = 2) 52,
  case 2:
                                                            51
                                                                 jmp 86
     (Teardown Router)
                                                                 (Teardown Router)
     \langle \text{Setup Router (NIC B)} \rangle
                                                                 \langle \text{Setup Router (NIC B)} \rangle
     break;
  default:
     . . .
     break;
}
```

Figure 4.7: Comparison of C code and Functional specification Router Manager control flow

In Subsection 4.1.4 and Subsection 4.1.5 we will look at the two main tasks from the Router Manager shown above, 'Teardown Router' and 'Setup Router'. In the subsections following those, we will go into more detail examining some selected instructions from the two tasks.

4.1.4 Example: Teardown Router

As mentioned in Subsection 4.1.3, each of the branches of the Router Manager program switch will first teardown a currently established Router (if there is one). Figure 4.8 gives the Router Manager Functional specification representation of that teardown process. The process of tearing down the Router is identical in the three branches, however the line numbers and *jmp* and *cnd_jmp* targets (not branch conditions) will be different, here α and β .

The Router teardown process is one loop nested inside another, the outer loop iterates over the various untyped sizes and the inner loop iterates over the non-free untyped capabilities. At each

iteration of the inner loop a non-free untyped capability is freed by revoking it (revocation is discussed in Section 2.2 and Subsection 4.1.6). At each iteration of the outer loop, *all* non-free capabilities of the current size are freed.

In addition to destroying all of the kernel objects and capabilities associated with the old Router, we will sanitise all network interfaces. This is done to ensure that no residual information from communications with one instantiation of the Router can bleed into a future instantiation. Santising one interface on the present architecture consists of restarting it.

Figure 4.8: Router Manager Functional specification Router teardown.

The teardown process in the Router Manager Functional specification (see Subsection 4.1.4) mostly consists of revoking untyped objects which are kept in a two-dimensional array *untyped_items*, where the first dimension is untyped object size, and the second is just the numbering untypeds. Revoking all the untyped items therefore requires iterating over both the sizes of untypeds and the number of untypeds of each size. At each iteration of the loop a kernel call is performed- meaning that the kernel must run- so the loop cannot be done in one Router Manager Functional specification intruction.

The instruction executed at each iteration of the inner loop is *revoke_ith_untyped*, explained in Subsection 4.1.6. The *sac_abstract_state* variable *var_i* contains the iterator for untyped size and *free_items* is used as the iterator for the inner loop. *free_items* is an array containing the number of free (as in, unallocated) untyped items for each untyped size, so, for a given size, we are iterating over each non-free untyped of that size.

4.1.5 Example: Setup Router

In Subsection 4.1.3 we discussed how two of the three branches of the overall Router Manager program flow will set up a Router thread. Figure 4.9 illustrates how the Router Manager Functional specification performs the setup. In the two 'setup' branches, the line numbers and *jmp* and *cnd_jmp* targets (not the branch conditions) will be different, here they are α , β , γ , and δ .

Additionally, the outside network interface (either Network A or Network B) will vary between the two setup branches, here I is used in place of either 'A' or 'B'.

The process of setting up a new Router involves creating a new Thread Control Block for the Router thread and then populating it with a *CSpace*, virtual address space, and an *IPC Buffer* (see Section 2.2). Additionally we need to copy some capabilities into the Router's *CSpace*. There are two loops in the setup, the first maps a number of 'scratch' frames into the Router's virtual address space and the second copies a number of untyped capabilities into the Router's *CSpace*.

```
k_instr create_router_cspace_object
    k_instr mutate_router_cspace_cap
    u_instr (modify (\lambda(tc, rm). (tc, rm(var_i:=0))))
    jmp \beta

    k_instr create_router_scratch_space

    k_instr map_router_scratch_frame
    k_instr move_router_scratch_frame
    u_instr (modify (\lambda(tc, rm). (tc, rm(val_i:=val_i rm + 1))))
\beta: cnd_jmp (\lambda(tc, rm). val_i rm < CLIENT_NUM_SCRATCH_FRAMES) \alpha
    k_instr create_router_ipc_object
    k_instr map_router_ipc_frame
    k_instr create_router_tcb_object
    k_instr setup_router_tcb
    k_instr set_router_tcb_ip_sp
    k_instr copy_time_t_client
    u_instr (modify (\lambda(tc, rm). (tc, rm(val_i:=0))))
    jmp \delta

γ: k_instr copy_untyped_to_client
    u_instr (modify (\lambda(tc, rm). (tc, rm(val_i:=val_i rm + 1))))
\delta: cnd_jmp (\lambda(tc, rm). val_i rm < CLIENT_NUM_UNTYPED_ITEMS) \gamma
    k_instr copy_asid_pool_to_client
    k_instr copy_nic_I_to_client
    k_instr copy_nic_D_to_client
    k_instr copy_tcb_cap_to_client
    k_instr move_ipc_buffer_cap_to_client
    k_instr mutate_cspace_to_client
```

Figure 4.9: Router Manager Functional specification Router setup.

4.1.6 In-depth *instr* examples

The previous sections of this document have laid out the infrastructure for expressing the Router Manager program in the Functional specification as well as the intended structure and purpose of the Router Manager program. This sub-section aims to bring that altogether and to provide the reader with some specific examples of code from the Router Manager program, how we wrote it and what it does in some detail. A line-by-line treatment would be very long so I have selected to present a few examples which should cover most 'sorts' of tasks in the program, and how they work.

The Router Manager Functional specification represents the Router Manager program as a list of instructions (unstructured code in the style of assembly), as discussed in Section 3.1. In this part of the document I will discuss several such 'instructions' from the Router Manager Functional specification's representation of the Router Manager program, *router_manager*. These particular examples also appear in Figure 4.8 and Figure 4.9. The full definition of *router_manager* can be found in Appendix A.

The very first thing that the Router Manager does, after booting, is to wait on the Control Network for a message (for more details see Section 2.4) which will contain either 'teardown Router', 'create Network A Router', and 'create Network B Router'. Here we will treat the process from the perspective of the Router Manager Functional specification. Figure 4.10 gives the definition in terms of a libseL4 interface function, *seL4_Wait*. The call must be lifted from *tc_monad* (see Subsection 4.2.1) to *rm_monad* since the libseL4 interface functions are defined across a generic thread state and not Router Manager specific state (see Subsection 3.2.2). *seL4_Wait* will be treated in depth in Subsection 4.2.4. The overall effect of this instruction is to set some registers up for the kernel, and then to trap into the kernel suspending itself until it receives a message.

```
definition
  wait_for_message ::
   (event option) rm_monad
  where
  wait_for_message =
   do command_endpoint' ← gets (λ(tc, rm). command_endpoint rm);
      lift_tc_rm (seL4_Wait command_endpoint' 0)
   od
```

Figure 4.10: Definition of *wait_for_message*

The main activity of *revoke_ith_untyped* is calling libseL4 interface function *seL4_CNode_Revoke*, via the Router Manager Functional specification abbreviation *revoke_untyped*, which lifts the libseL4 function and supplies some constant arguments.

revoke_untyped is given in Figure 4.12. Calling *seL4_CNode_Revoke* should set some registers up and trap into the kernel, so that the kernel can revoke the specified untyped object- destroying any kernel objects residing in the untyped object along with any capabilities referring to them. In the context of the whole program; when this is done on untyped memory objects which were used to set up a previous Router thread, it will result in stripping of the Router of its capabilities and then shut the Router thread itself down.

```
definition
  revoke_ith_untyped ::
  (event option) rm_monad
  where
  revoke_ith_untyped =
   do i \leftarrow gets (\lambda(tc, rm). var_i rm);
      free_items' \leftarrow gets (\lambda(tc, rm). free_items rm);
      untyped_items' \leftarrow gets (\lambda(tc, rm). untyped_items rm);
      revoke_untyped RM_SelfCSpace
         ((untyped_items'!i)!(unat (free_items'!i)))
   od
                Figure 4.11: Definition of revoke_ith_untyped
definition revoke_untyped ::
  seL4\_CPtr \Rightarrow seL4\_CPtr \Rightarrow (event option) rm_monad
  where
  revoke_untyped dest_cspace dest_slot =
     lift_tc_rm (seL4_CNode_Revoke dest_cspace dest_slot
                                        (of_nat seL4_WordBits))
```

Figure 4.12: Definition of *revoke_untyped*.

The process of setting up a new Router is much more complex than tearing one down (see Subsection 4.1.5), so all of the tasks involved can't be broken down in this section. The example presented is the instruction that will create a kernel object for the Router's Thread Control Block (TCB) in the kernel. A TCB in the seL4 microkernel is quite close to what a reader will be familiar with from a Unix or other L4 TCB. Importantly it contains an Inter-Process Communication buffer (IPC buffer), the root CNode of the thread's CSpace (see Section 2.6), and the root of the thread's virtual address space. This important information will be populated in subsequent instructions-the create instruction should just provide us with a fresh TCB object, along with a capability to it.

Figure 4.13: Definition of *create_router_tcb_object*.

As shown in Figure 4.13, we first call a function to 'allocate' a small kernel object (see Subsection 4.1.7) by updating some internal bookkeeping on untyped memory. Next, we call the libseL4 interface function *seL4_Untyped_Retype*, via Router Manager Functional specification abbreviation *create_object. create_object*, as shown in Figure 4.14, is merely a wrapper which lifts the libseL4 interface function *seL4_Untyped_Retypr* to *rm_monad*.

4.1.7 Internal helper function examples

For the most part we tried to mirror the internal structure of the Router Manager C implementation for the Functional specification program. At some points this meant clumping user transitions together with a kernel call as shown in Subsection 4.1.6. When this happened, there were often internal helper functions in the C implementation so instead of 'inlining' these purely internal transitions they were represented as *rm_monad* computations to be used in other instructions.

get_message, shown in Figure 4.15, calls libseL4 helper function *seL4_GetMR* to retrieve some information from a message register. Message register being a user view abstraction of the Router Manager thread's IPC buffer. The structure of libseL4 and the functions it exposes is discussed in much more detail in section Section 4.2. In this case the message refered to is from the Control Network, and contains direction on what to do next.

Figure 4.15: Definition of *get_message*.

allocate_small_item is shown in Figure 4.16. This helper function is called to update bookkeeping information on untyped memory (see Subsection 4.1.6 for an example of its use). It is an abbreviation for *allocate_item*.

```
definition allocate_small_item ::
    seL4_CPtr rm_monad
    where
    allocate_small_item = allocate_item SIZE_4K
        Figure 4.16: Definition of allocate_small_item.
```

Figure 4.17 gives the definition for *allocate_item*. First, we decrease the number of 'free' untyped items of the given size in *free_items*. Second we find the next 'free' untyped item in *untyped_items* and return a pointer to it.

```
definition
```

```
allocate_item ::

untyped_size \Rightarrow seL4_CPtr rm_monad

where

allocate_item sz =

do modify (\lambda(tc, rm). (tc, free_items_update (\lambdax. x[sz := (x!sz) - 1]) rm));

ptr \leftarrow gets (\lambda(tc, rm).

((untyped_items rm)!sz)!(unat ((free_items rm)!sz)));

return ptr

od
```

Figure 4.17: Definition of *allocate_item*.

4.2 libseL4 specification

We have now defined the foundations of the Router Manager Functional specification, and shown some examples of how the Router Manager component is expressed up to where the kernel is called. Beyond that, there are libseL4 functions setting up and trapping into the kernel. To completely specify the Router Manager we need to specify the subset of libseL4 that the Router Manager uses.

Just as the Router Manager instructions are treated as indivisible monadic computations, so are the libseL4 functions. In the future there may be some reason to implement pre-emptible libseL4 interface functions, however it seems like a reasonable assumption that seL4 user systems (especially those that you would like to verify!) can separately be proven to have interference freedom from the rest of the system, with regards to their *abstract_thread_context* (see Subsection 4.2.2). If the *abstract_thread_context* is interference free then there is no harm in representing the libseL4

interface functions as atomic.

As with the previous section, it would be infeasible to present the entire specification of the libseL4 library here. The foundational definitions are given in this section along with examples which are meant to be illustrative of how the specification works.

4.2.1 The *tc_monad*

```
types 'a tc_monad = (thread_abstract_context, 'a) nondet_monad
Figure 4.18: Definition of tc_monad.
```

As the libseL4 interface functions are meant to be user program generic, they are defined across a thread state rather than the Router Manager, or Secure Access Controller, state. These computations are represented by the *tc_monad*, illustrated in Figure 4.18. Similarly to the Router Manager instructions, functions which trap into the kernel are expected to return some *SyscallEvent sy* of the appropriate sort, where the C implementation of that same function would have set up a software interrupt.

4.2.2 The *thread_abstract_context*

types thread_abstract_context = seL4_IPCBuffer × user_context
Figure 4.19: Definition of thread_abstract_context.

thread_abstract_context represents user state which is common to all seL4 user programs, and which is also necessary for the libseL4 library functions. Shown in Figure 4.19, the *seL4_IPCBuffer* component of the state is the Inter-Process Communication Buffer (IPC Buffer) of the user program thread and the *user_context* represents the user view of machine registers. The IPC Buffer is a buffer for messages sent to the user thread by the kernel or other processes. The *user_context* type comes from the seL4 kernel Abstract specification. Defining the libseL4 functions in terms of *user_context* is convenient because we wish to use the seL4 kernel Abstract specification in our implementation of the overall system process.

4.2.3 libseL4 specification structure

The seL4 microkernel provides six syscalls and several dozen capability invocations. Capability invocations are particular instances of syscalls: the user gives the kernel some capability while making a syscall and the kernel interprets that as some *invocation* on the given capability. For

example, *seL4_Wait* is a syscall and *seL4_CNode_Revoke* is a capability invocation (see Subsection 4.1.6). The libseL4 library provides C functions that will set up both the syscalls and the capability invocations. Because the capability invocations depend on the syscalls, we can view the libseL4 syscall setup functions as the libseL4 interface to the kernel and the capability invocations as the user's interface to libseL4. Because of this, in the libseL4 specification the syscalls are part of the *seL4 interface* while the capability invocations are part of the *libseL4 interface*.

The C implementation of the libseL4 library can be broken down conceptually into a seL4 interface, a libseL4 interface, and a set of types which are used to represent the user view of kernel objects as well as to pass data into the kernel.

We implemented the libseL4 specification as a package of types that user programs can use to represent seL4 kernel objects, on top of which are functions to set and get data from the *thread_abstract_context* (see Subsection 4.2.5), seL4 interface functions (syscalls), and then libseL4 interface functions.

Only the parts of the libseL4 interface used by the Router Manager program have been implemented, however what has been implemented is generic enough that other user programs could easily use the same libseL4 specification, assuming architecture compatibility. Filling in the rest of the interface would be straight-forward and not time-consuming, since everything that the functions depend on has been set up in the process of getting a subset working.

4.2.4 In-depth libseL4 examples

Presented here are some examples of libseL4 and seL4 interface functions. Their purpose may not be immediately obvious from the definitions: in each case some information is marshalled and loaded into registers where the kernel can find it and special values are loaded into other registers so that the kernel can determine the user's intent. Then a software interrupt is generated, trapping into the kernel so that the desired operation can be performed.

Careful readers will notice that the libseL4 interface functions load data into message registers which the seL4 interface functions then put into machine registers. This is definitely overly complex and costly to performance but exists in the current implementation for historical reasons, and because code is automatically generated. Correcting this behaviour is very much outside of the scope of this thesis project, however it was noticed and modifying these definitions when or if this behaviour is changed would not be burdensome (indeed it would make the definitions smaller).

seL4_Call, shown in Figure 4.20, is part of what we call the seL4 interface of libseL4, it provides a convenient way to generate a seL4 syscall. We take some arguments and load them into registers, then move some data from message registers into machine registers (or at least the user's view

of the machine registers).

```
definition seL4_Call ::
  seL4\_CPtr \Rightarrow seL4\_MessageInfo \Rightarrow (event option) tc_monad
  where
  seL4_Call dest msgInfo =
      do modify (\lambda(i,u). (i, u (RO := dest)));
         modify (\lambda(i,u). (i, u (R1 := msgInfo)));
         MR0 \leftarrow gets (seL4_GetMR 0);
         modify (\lambda(i,u). (i, u (R2 := MR0)));
         MR1 \leftarrow gets (seL4_GetMR 1);
         modify (\lambda(i,u). (i, u (R3 := MR1)));
         MR2 \leftarrow gets (seL4\_GetMR 2);
         modify (\lambda(i,u). (i, u (R4 := MR2)));
         MR3 \leftarrow gets (seL4_GetMR 3);
         modify (\lambda(i,u). (i, u (R5 := MR3)));
         modify (\lambda(i,u). (i, u (R7 := seL4_SysCall)));
         return (Some (SyscallEvent SysCall))
     od
```

Figure 4.20: Definition of *seL4_Call*.

Figure 4.21 presents *seL4_Wait*, which is also a seL4 interface function. It sets up a seL4 syscall similarly to *seL4_Wait*. It is often used directly by user programs, for example the Router Manager (see Figure 4.21).

Figure 4.21: Definition of *seL4_Wait*.

seL4_CNode_Revoke is an example of a capability invocation or *libseL4 interface* function. Shown in Figure 4.22, its purpose is to revoke any children of the invoked capability. In the case of a capability to an untyped memory object, any kernel object created in that untyped will be a child and so revoking the untyped capability will destroy any objects created inside. The Router Manager program uses this function in tearing down the Router thread, by repeatedly revoking untyped capabilities (see Figure 4.8).

The purpose of *seL4_Untyped_Retype*, shown in Figure 4.23 is to create a new kernel object in the memory covered by invoked untyped capability. In practice it places some data into message

```
definition
  seL4_CNode_Revoke ::
  seL4_CNode ⇒ seL4_Word ⇒ 8 word
  ⇒ (event option) tc_monad
  where
  seL4_CNode_Revoke service index depth =
    do seL4_SetMR 0 index;
      seL4_SetMR 1 (of_nat (unat depth) && 0xff);
      seL4_Call service (seL4_MessageInfo 2 0 (of_nat (unat CNodeRevoke)) 0)
    od
```

Figure 4.22: Definition of seL4_CNode_Revoke.

registers and calls *seL4_Call*, which will do more setup and generate a software interrupt. This function is used extensively in the Router Manager's setup of a Router thread (see Figure 4.9). *seL4_Untyped_Retype* is also a *libseL4 interface* function.

definition

Figure 4.23: Definition of *seL4_Untyped_Retype*.

4.2.5 libseL4 helper function examples

The functions presented in this sub-section are used heavily in the seL4 and libseL4 interface functions of the libseL4 specification. They are quite straight-forward but are presented so that the continuing examples from Subsection 4.1.6 to Subsection 4.2.4 are completely described (the complete definitions of all other parts of the specifications are given in Appendix A, Appendix B, and Appendix C).

The function *seL4_GetMR* is heavily used in the definition of *seL4_Call*, is shown in Figure 4.24. It simply returns the value inside some message register in the *thread_abstract_context* of the user program thread. It does not return any event because it is used inside the 'atomic' libseL4 and seL4 interface functions.

```
definition seL4_GetMR ::

nat \Rightarrow thread_abstract_context \Rightarrow seL4_Word

where

seL4_GetMR n = (\lambda(i, u). (ipc_msg i) ! n)

Figure 4.24: Definition of seL4_GetMR.
```

The function *seL4_SetMR*, shown in Figure 4.25, will set the given message register to the given value. It does not return any event because it too is used inside of other 'atomic' functions.

Figure 4.25: Definition of *seL4_SetMR*.

4.2.6 Lifting from *tc_monad* to *rm_monad*

The aim of the libseL4 specification is to have a generic interface that any user program specification can use. The libseL4 functions, however, are defined in terms of the *tc_monad*, and it is unlikely that user program specifications could be fully captured by *tc_monad*.

The Router Manager functional specification, for example, uses *sac_abstract_state* in its *rm_monad*. Therefore in order for Router Manager program specification to use the libseL4 functions we need a way to lift functions from the *tc_monad* to *rm_monad*.

The function *lift_tc_rm*, shown in Figure 4.26, lifts a computation from the *tc_monad* to the *rm_monad*. It carries the *tc_monad* computation on the *thread_abstract_context* part of the Router Manager state, and then recombines the resulting state with the rest of the Router Manager state and returns the return value of the *tc_monad* computation.

```
definition

lift_tc_rm ::

'a tc_monad \Rightarrow 'a rm_monad

where

lift_tc_rm f =

(\lambda(tc, rm). ((\lambda(e, tc'). (e, (tc', rm))) \text{ 'fst (f tc), snd (f tc))})

Figure 4.26: Definition of lift_tc_rm.
```

4.3 System process

As mentioned in Chapter 1, showing that the Router Manager Functional specification is a refinement of the Router Manager Security model will require some way to create a process from each. We did this by defining a generic system automaton in terms of a user specification and a kernel specification. The idea of the automaton is to capture all of the potential executions of the user and kernel with arbitrary scheduling, of course covering all user to kernel transitions and kernel to user transitions.

Aside from the main goal of this thesis project- to formally specify the Router Manager- we also implemented the generic system automaton to allow any system processes to be defined and we instantiated it with a system process for the Router Manager Functional specification. This section presents the generic system automaton, the system process as it is defined for the Router Manager Functional specification, how the kernel specification is defined, and how the user specification is defined.

4.3.1 The global_automaton

In this sub-section we will introduce some necessary types and then the definition of the *global_automaton* itself.

The *mode* datatype, depicted in Figure 4.27, represents our notion of which mode the entire system is running in at a given time. At all times the kernel is considered to be in one of: *Idle mode, Kernel mode,* or *User mode*. The distinction is rather self-explanatory, the system is executing in *Idle mode* when neither a thread nor the kernel is running, *Kernel mode* when the kernel is running (it cannot be interrupted by the scheduler or pre-empted), and *User mode* when some user thread is running. This is illustrated in Figure 3.3.

datatype mode = UserMode | KernelMode | IdleMode Figure 4.27: Definition of *mode*.

As shown in Figure 4.28, user_mem represents a user view of system memory. 32 bit addresses

potentially map to data in memory, dependant on whether the current user thread should be able to access that address or not. Along with *user_context* and the abstract *Instruction pointer*, this makes up the user view of the machine.

types user_mem = word32 \Rightarrow word8 option Figure 4.28: Definition of *user_mem*.

The type *user_context* is meant to capture the user program's view of the system's registers. As can be seen in Figure 4.29, the type is a map from *register* to *data*. The types of *register* and *data* are not given here for the sake of length, however the former is just an abstract type for a register's name (for example *R0*) and the latter is a 32 bit word.

types user_context = register \Rightarrow data Figure 4.29: Definition of user_context.

The *global_automaton* is intended to, given a user specification and kernel specification, return all of the system's transitions. For consistency with the refinement machinery from the seL4 refinement, the transitions are broken up structurally into the different *global_transition* types. A *global_transition* is either a *KernelTransition*, *UserTransition*, *UserEventTransition*, or an *IdleEventTransition*, however the *UserEventTransition* is no longer used because of the user having its own specification in this definition.

As shown in Figure 4.30 the definition is broken up into different types of transitions, each giving the set of possible system transitions. For example, the kernel transitions consist of the union of some transitions from *KernelMode* to *UserMode* and from *KernelMode* to *IdleMode*. This is because the kernel is expected to not generate kernel events and hence once the system enters *KernelMode* it can only transition to *UserMode* or *IdleMode* (presumably if there are no user threads ready to run).

For the definition of 'allowed' kernel transitions, *global_automaton* appeals to the given kernel specification *kernel_call*, and for 'allowed' user transitions it appeals to the user specification *user_transition*. Additionally we allow the user threads to generate an *Interrupt* event at any time, since these will not be explicitly returned by the user specification. Finally, the only 'allowed' *IdleMode* transition is to *KernelMode* via an *Interrupt* event.

The type of *global_automaton* is not concrete, there is a type variable 'k which should be instantiated with a seL4 kernel state type, for example the seL4 Abstract specification *state*. The types of *user_transition* and *kernel_call* also depend on 'k, so they must use the same kernel state. Fully applied, *global_automaton* is a relation on 'k *global_state*, representing the possible transitions between 'k *global_state* states.

```
global_automaton ::
  ((user_context \times word32 \times k) \times (user_context \times word32 \times k)
    \times mode \times event option) set
   \Rightarrow (event \Rightarrow ((user_context \times 'k) \times mode \times (user_context \times 'k)) set)
   \Rightarrow global_transition \Rightarrow ('k global_state \times 'k global_state) set
where
global_automaton user_transition kernel_call t \equiv case t of
  KernelTransition \Rightarrow
     { ( ((uc, ip, s), KernelMode, Some e),
          ((uc', uc' LR_svc, s'), UserMode, None) ) |uc ip s uc' s' e.
            ((uc, s), UserMode, (uc', s')) \in kernel_call e }
   \cup { ((uc, ip, s), KernelMode, Some e),
          ((uc', ip', s'), IdleMode, None) ) |uc ip s uc' ip' s' e.
            ((uc, s), IdleMode, (uc', s')) \in kernel_call e \}
| UserTransition \Rightarrow
     { ( (g, UserMode, None), (g', UserMode, e) ) \mid g g' e.
     (g, g', UserMode, e) ∈ user_transition}
   \cup { ((uc, ip, g), UserMode, None),
          ((uc'(FaultInstruction := ip'), ip', g'), KernelMode, Some ev) )
          luc ip g uc' ip' g' ev.
            ((uc, ip, g), (uc', ip', g'), KernelMode, Some ev) \in
              user_transition}
   ∪ { ( (g, UserMode, None), (g, KernelMode, Some Interrupt) ) |g. True }
| IdleEventTransition \Rightarrow
   { ( (g, IdleMode, None), (g, KernelMode, Some Interrupt) ) |g. True }
```

Figure 4.30: Definition of *global_automaton*.

4.3.2 The *ADT_I*

Given the generic definitions that we implemented, as presented in Subsection 4.3.1, we may now define the entire system process for the seL4 microkernel and the Router Manager program. Shown in Figure 4.31, the system process is called *ADT_I*. *ADT_I* is a record containing three functions: *Init*, *Fin*, and *Step*.

Init is a function which initializes the system process by setting up a *state global_state* which corresponds to a given *observable* state (see Figure 2.2). For any system process, we must choose some kernel view of the system, this will be the state of the kernel specification. *observable* state is a representation of the system state which is not dependant on that kernel view of the system. Because it is not dependant on any internal kernel view of the system, the *observable* state is comparable between system processes which could have different kernel specifications. As mentioned in Subsection 4.3.1, the *global_automaton* is a relation over some 'k global_state, which in the case of the *ADT_I* becomes concrete type *state global_state*.

The function *Fin* does more or less the reverse of *Init*, it finalizes the system process. *Fin* will give the *observable* state which corresponds to a given *state global_state*. In forward simulation, for a particular system process, we will use *Init* and *Fin* to encode an arbitrary state into that process's internal representation. Then after taking a step, we decode the resulting states into something comparable with another system process.

The *Init* and *Fin* functions given in *ADT_I*, *Fin_A* and *Init_A*, are both taken from the seL4 Abstract specification automaton that our definition has been adapted from, since *ADT_I* uses the Abstract kernel specification.

The third field, *Step*, is probably the most interesting. *Step* encodes all of the valid transitions of the system process. For each *global_transition*, it gives a set of all system transitions of that type. For *ADT_I*, *Step* is *global_automaton* instantiated in terms of *do_user_op*, *uop_I*, and *kernel_call_A*.

The function *uop_I* is the unlifted user specification and is discussed in Subsection 4.3.4. *do_user_op* is a helper definition from the generic automaton machinery I adapted, it lifts *uop_I* to transitions on *state global_state*. *kernel_call_A* is the kernel specification for the seL4 Abstract specification, and is discussed in Subsection 4.3.3.

The type of *ADT_I* uses abbreviation *data_type* for historical reasons. For type variables (α , β , γ), *data_type* will be a record with fields *Init* :: $\beta \Rightarrow \alpha$, *Fin* :: $\alpha \Rightarrow \beta$, and *Step* :: $\gamma \Rightarrow (\alpha \times \alpha)$.

```
definition
   ADT_I :: (state global_state, observable, global_transition) data_type
where
   ADT_I =
    ([ Init = λs. Init_A,
        Fin = Fin_A,
        Step = global_automaton (do_user_op uop_I) kernel_call_A])
        Figure 4.31: Definition of ADT_I.
```

4.3.3 Kernel specification

The definition of *kernel_call_A* is given here, in Figure 4.32, for the curious. It gives the type of our seL4 Abstract kernel specification. The definition obviously does not define the behaviour of the kernel directly- there are many thousands of lines of definition and proof to it.

4.3.4 User specification

The user specification for our system process in general just needs to specify the valid *user_state* transitions and which sort of *event* they generate. We say *user_state* transitions rather than *state*

```
definition
   kernel_call_A ::
   event ⇒ ((user_context × state) × mode × (user_context × state)) set
   where
   kernel_call_A e ≡
      {(s, m, s'). s' ∈ fst (split (kernel_entry e) s) ∧
            m = (if ct_running (snd s') then UserMode else IdleMode)}
```

Figure 4.32: Definition of *kernel_call_A*.

global_state transitions, since there is a generic function, *do_user_op*, to lift to *state global_state*.

In the particular case of the Router Manager Functional specification, we have already examined the *sac_abstract_state* and the *rm_monad* in Section 4.1. What is presented in this subsection is how we lift *rm_monad* computations to *user_state* transitions and how we lift *user_state* transitions to *state global_state* transitions.

uop_I, as shown in Figure 4.33, decodes the program text of the currently running thread from user memory giving a list of *user_state* transitions. It then decodes the program counter from the user state and 'executes' that instruction from the list, at the current user state, returning the result.

definition

```
uop_I ::
user_state ⇒ (user_state × event option) set
where
uop_I us =
  ((decode_program_text
      (project_program_text us))!(project_program_counter us)) us
```

Figure 4.33: Definition of *uop_I*.

decode_program_text maps a region of data, containing machine instructions, to a list of *user_state* instructions. In practice the result will be an arbitrary list for *untrusted* processes or the *router_manager* program, lifted to *user_state*, for the *trusted* Router Manager process. Therefore when reasoning about the system process for the Router Manager Functional specification, transitions when an *unstrusted* process is running can be anything and transitions when the Router Manager process is running are those from Functional specification.

project_program_text projects the region of data, containing the currently running process's program text, out of the user process's entire view of memory (contained in *user_state*).

project_program_counter projects the program counter out of the *user_state*, which will be the abstract *instruction_pointer*.

Chapter 5

Conclusions

The goal of this project was to formally specify the functionality of the Router Manager component of the Secure Access Controller, so that it could serve as a midway point in the refinement from the SAC Security model to its C implementation. Since the C implementation uses functions from the common seL4 user library libseL4, it was necessary along the way to specify the behaviour of a libseL4 subset as well.

5.1 Contributions

We succeeded in implementing the Router Manager Functional specification, leaving the specification of libseL4 as a separate object. Additionally, we implemented the subset of libseL4 needed for the Router Manager program as well as the framework for showing refinement. The libseL4 subset was designed to be general enough that another seL4 user program could use its functions, by writing a simple lift from the libseL4 state to the user program's state. The refinement framework was designed so that, similarly, it could be used to express refinement between other seL4 user program specifications.

For the above reasons, this project has been quite helpful as a case study in seL4 user program verification.

5.2 Future Work

The next steps after this thesis will be to define a system process for the SAC Security model Router Manager program and then to prove that the Router Manager Functional specification refines it. This will most likely necessitate modifying the SAC Security model. Security model transitions need to be changed so that they can be lifted to *user_state* transitions which return an *event option*.

Further into the future, we will need to show that the C implementation of the Router Manager refines the Router Manager Functional specification. Doing so will likely require adapting the system process definitions a little. Because where the system process for the Functional specification takes the seL4 Abstract specification, the system process for the Router Manager's C implementation should take a different kernel specification.

Appendix A

Router Manager Functional Specification

theory RouterManager_I

imports libseL4Interface

begin

A.1 The Router Manager Functional specification program

Router Manager constant definitions

```
definition CLIENT_CPTR_TCB :: "nat" where "CLIENT_CPTR_TCB = 1"
definition CLIENT_CPTR_CSPACE :: "nat" where "CLIENT_CPTR_CSPACE = 2"
definition CLIENT_CPTR_IPC_BUFFER :: "nat" where "CLIENT_CPTR_IPC_BUFFER = 4"
definition CLIENT_CPTR_ASID_POOL :: "nat" where "CLIENT_CPTR_ASID_POOL = 6"
definition CLIENT_CPTR_IO_PORT :: "nat" where "CLIENT_CPTR_IO_PORT = 7"
definition CLIENT_CPTR_TIMER_ENDPOINT :: "nat" where "CLIENT_CPTR_TIMER_ENDPOINT = 8"
definition CLIENT_CPTR_FIRST_UNTYPED_ITEM :: "nat" where
    "CLIENT_CPTR_FIRST_UNTYPED_ITEM = 1024"
definition CLIENT_CPTR_FIRST_SCRATCH_ITEM :: "nat" where
    "CLIENT_CPTR_FIRST_SCRATCH_ITEM = 2048"
```

definition UNTYPED_ITEM_SIZE_BITS :: "nat" where "UNTYPED_ITEM_SIZE_BITS = 20"

```
definition CLIENT_NUM_SCRATCH_FRAMES :: "nat" where "CLIENT_NUM_SCRATCH_FRAMES = 32"
definition CLIENT_NUM_UNTYPED_ITEMS :: "nat" where "CLIENT_NUM_UNTYPED_ITEMS = 100"
definition CLIENT_CSPACE_SIZE_BITS :: "nat" where
```

"CLIENT_CSPACE_SIZE_BITS = (UNTYPED_ITEM_SIZE_BITS - seL4_SlotBits)" definition CLIENT_IPC_BUFFER_VADDR :: "nat" where

"CLIENT_IPC_BUFFER_VADDR = 0xd0002000"

definition CLIENT_SCRATCH_VADDR :: "nat" where "CLIENT_SCRATCH_VADDR = 0xd0003000"
definition CLIENT_ENTRY_POINT :: "nat" where "CLIENT_ENTRY_POINT = 0x00010000"
definition CLIENT_PRIORITY :: "nat" where "CLIENT_PRIORITY = 128"

```
definition RM_Null :: "seL4_Word" where "RM_Null = 0"
definition RM_SelfCSpace :: "seL4_Word" where "RM_SelfCSpace = 2"
definition RM_AsidPool :: "seL4_Word" where "RM_AsidPool = 6"
definition RM_IOPort :: "seL4_Word" where "RM_IOPort = 7"
```

```
definition RM_Client_Tcb :: "seL4_Word" where "RM_Client_Tcb = 1000"
definition RM_Client_CSpace :: "seL4_Word" where "RM_Client_CSpace = 1001"
definition RM_Client_IpcBuffer :: "seL4_Word" where "RM_Client_IpcBuffer = 1002"
definition RM_TempSlot :: "seL4_Word" where "RM_TempSlot = 1003"
```

```
types untyped_size = nat
definition SIZE_4K :: "nat" where "SIZE_4K = 0"
definition SIZE_1MB :: "nat" where "SIZE_1MB = 1"
definition NUM_SIZES :: "nat" where "NUM_SIZES = 2"
```

A.2 The sac_abstract_state

```
record sac_abstract_state =
```

ip :: seL4_Word

```
untyped_items :: "seL4_CPtr list list"
free_items :: "seL4_Word list"
total_items :: "seL4_Word list"
info :: seL4_MessageInfo
message :: seL4_Word
var_i :: nat
RM_RWGC_NIC_A_cap :: seL4_CPtr
RM_RWGC_NIC_D_cap :: seL4_CPtr
client_vspace :: seL4_CPtr
virtualirq_endpoint :: seL4_CPtr
```

command_endpoint :: seL4_CPtr

A.3 The rm_monad and Router Manager execution

```
types 'a rm_monad = "((thread_abstract_context × sac_abstract_state), 'a) nondet_monad"
types instr = "(event option) rm_monad"
types program = "instr list"
```

```
lift_tc_rm ::

"'a tc_monad \Rightarrow 'a rm_monad"

where

"lift_tc_rm f =

(\lambda(tc, rm). ((\lambda(e, tc'). (e, (tc', rm))) ' fst (f tc), snd (f tc)))"
```

```
definition incr_pc ::
    "unit rm_monad"
    where
    "incr_pc = modify (λ(tc, rm). (tc, ip_update (λip. ip + 1) rm))"

definition k_instr ::
    "instr ⇒ instr"
    where
    "k_instr = id"
```

```
definition u_instr ::
  "'a rm_monad \Rightarrow instr"
  where
  "u_instr m =
     do m;
         incr_pc;
         return None
     od"
definition jmp ::
  " seL4_Word \Rightarrow instr"
  where
  "jmp n =
     do modify (\lambda(tc, rm). (tc, ip_update (\lambdax. x + 1) rm));
         return None
     od"
definition cnd_jmp ::
  "((thread_abstract_context \times sac_abstract_state) \Rightarrow bool) \Rightarrow seL4_Word \Rightarrow instr"
  where
  "cnd_jmp P n =
      (\lambdas. (if (P s)
             then
               jmp n
             else return None) s)"
```

A.4 Router Manager helper functions

```
definition
  allocate_item ::
    "untyped_size ⇒ seL4_CPtr rm_monad"
    where
    "allocate_item sz =
        do modify (λ(tc, rm). (tc, free_items_update (λx. x[sz := (x!sz) - 1]) rm));
        ptr ← gets (λ(tc, rm). ((untyped_items rm)!sz)!(unat ((free_items rm)!sz)));
        return ptr
        od"
    definition allocate_large_item ::
```

```
"seL4_CPtr rm_monad"
```

```
where
  "allocate_large_item = allocate_item SIZE_1MB"
definition allocate_small_item ::
    "seL4_CPtr rm_monad"
   where
    "allocate_small_item = allocate_item SIZE_4K"
```

```
get_message ::
"unit rm_monad"
where
"get_message =
  do m ← gets (λ(tc, rm). seL4_GetMR 0 tc);
    modify (λ(tc, rm). (tc, rm(|message := m|)));
    return ()
  od"
```

A.5 Router Manager instructions

```
definition get_cspace_capdata ::
   "seL4_Word \Rightarrow seL4_CapData"
  where
   "get_cspace_capdata cspace_size_bits =
       seL4_CapData (of_nat (seL4_WordBits - unat cspace_size_bits)) 0"
definition try_copy_cap ::
   \texttt{"seL4\_CPtr} \Rightarrow \texttt{seL4\_CPtr} \Rightarrow \texttt{seL4\_CPtr} \Rightarrow \texttt{seL4\_CPtr} \Rightarrow \texttt{seL4\_CPtr} \Rightarrow \texttt{seL4\_CPtr} \Rightarrow \texttt{seL4\_CPtr}
    \Rightarrow (event option) rm_monad"
  where
   "try_copy_cap dest_cspace dest_slot src_cspace src_slot rights =
       lift_tc_rm (seL4_CNode_Copy dest_cspace dest_slot (of_nat seL4_WordBits)
                                                  src_cspace src_slot (of_nat seL4_WordBits) rights)"
definition copy_cap ::
   \texttt{"seL4\_CPtr} \ \Rightarrow \ \texttt{seL4\_CPtr} \ \Rightarrow \ \texttt{seL4\_CPtr}
    \Rightarrow (event option) rm_monad"
  where
   "copy_cap dest_cspace dest_slot src_cspace src_slot rights =
       try_copy_cap dest_cspace dest_slot src_cspace src_slot rights"
```

```
definition create_object ::
  \texttt{"seL4\_CPtr} \ \Rightarrow \ \texttt{seL4\_CPtr} \ \Rightarrow \ \texttt{seL4\_CPtr} \ \Rightarrow \ \texttt{seL4\_Word} \ \Rightarrow \ \texttt{seL4\_Word} \ \Rightarrow \ \texttt{seL4\_Word}
   \Rightarrow (event option) rm_monad"
  where
  "create_object untyped_item dest_cnode dest_slot targettype targetsize =
      lift_tc_rm (seL4_Untyped_Retype untyped_item targettype targetsize dest_cnode
                        0 0 dest_slot 1)"
definition map_frame ::
  \texttt{"seL4\_CPtr} \ \Rightarrow \ \texttt{seL4\_CPtr} \ \Rightarrow \ \texttt{seL4\_CPtr} \ \Rightarrow \ \texttt{seL4\_CPtr} \ \Rightarrow \ \texttt{seL4\_CPtr}
   \Rightarrow (event option) rm_monad"
  where
  "map_frame vspace frame vaddr perms =
      lift_tc_rm (seL4_ARM_Page_Map frame vspace vaddr
                        perms seL4_ARM_Default_VMAttributes)"
definition move_cap ::
  "seL4\_CPtr \Rightarrow seL4\_Word \Rightarrow seL4\_CPtr \Rightarrow seL4\_CPtr \Rightarrow (event option) rm_monad"
  where
  "move_cap dest_cspace dest_slot src_cspace src_slot =
      lift_tc_rm (seL4_CNode_Move dest_cspace dest_slot (of_nat seL4_WordBits)
                                      src_cspace src_slot (of_nat seL4_WordBits))"
definition revoke_untyped ::
  "seL4_CPtr \Rightarrow seL4_CPtr \Rightarrow (event option) rm_monad"
  where
  "revoke_untyped dest_cspace dest_slot =
      lift_tc_rm (seL4_CNode_Revoke dest_cspace dest_slot
                                             (of_nat seL4_WordBits))"
definition mutate_cap ::
  "seL4\_CPtr \Rightarrow seL4\_Word \Rightarrow seL4\_CPtr \Rightarrow seL4\_Word \Rightarrow seL4\_CapData
   \Rightarrow (event option) rm_monad"
  where
  "mutate_cap dest_cspace dest_slot src_cspace src_slot data =
      lift_tc_rm (seL4_CNode_Mutate dest_cspace dest_slot
                                             (of_nat seL4_WordBits) src_cspace src_slot
                                             (of_nat seL4_WordBits) data)"
```

```
setup_tcb ::
"seL4_CPtr ⇒ seL4_CPtr ⇒ seL4_CapData ⇒ seL4_CPtr ⇒ seL4_Word ⇒ seL4_Word
⇒ seL4_CPtr
⇒ (event option) rm_monad"
where
"setup_tcb tcb cspace cspace_data vspace priority buffer_vaddr buffer_frame =
lift_tc_rm (seL4_TCB_Configure tcb 0 (of_nat (unat priority)) cspace
cspace_data vspace 0
buffer_vaddr buffer_frame)"
```

definition

definition

```
sanitise_network_cards ::
"(event option) rm_monad"
where
"sanitise_network_cards = undefined"
```

definition

```
copy_nic_A_to_client ::
"(event option) rm_monad"
where
"copy_nic_A_to_client =
do cap ← gets (λ(tc, rm). RM_RWGC_NIC_A_cap rm);
try_copy_cap RM_Client_CSpace cap RM_SelfCSpace
cap seL4_AllRights
```

```
od"
```

```
copy_nic_B_to_client ::
```

```
"(event option) rm_monad"
where
"copy_nic_B_to_client =
  do cap ← gets (λ(tc, rm). RM_RWGC_NIC_B_cap rm);
    try_copy_cap RM_Client_CSpace cap RM_SelfCSpace
        cap seL4_AllRights
  od"
```

```
copy_nic_D_to_client ::
"(event option) rm_monad"
where
"copy_nic_D_to_client =
  do cap ← gets (λ(tc, rm). RM_RWGC_NIC_D_cap rm);
    try_copy_cap RM_Client_CSpace cap RM_SelfCSpace
        cap seL4_AllRights
  od"
```

definition

od"

definition

```
create_router_ipc_object ::
"(event option) rm_monad"
```

```
where
```

definition

od"

definition

```
setup_router_tcb ::
"(event option) rm_monad"
where
"setup_router_tcb =
  do vs ← gets (λ(tc, rm). client_vspace rm);
    setup_tcb RM_Client_Tcb RM_Client_CSpace
        (get_cspace_capdata (of_nat CLIENT_CSPACE_SIZE_BITS)) vs
        (of_nat CLIENT_PRIORITY) (of_nat CLIENT_IPC_BUFFER_VADDR)
        RM_Client_IpcBuffer
    od"
```

definition create_router_scratch_space ::

```
"(event option) rm_monad"
  where
  "create_router_scratch_space =
   create_object item RM_SelfCSpace RM_TempSlot
                    seL4_NonArchObjectTypeCount 0
   od"
definition map_router_scratch_frame ::
  "(event option) rm_monad"
  where
  "map_router_scratch_frame =
   do vs \leftarrow gets (\lambda(tc, rm). client_vspace rm);
      i \leftarrow gets (\lambda(tc, rm). var_i rm);
      map_frame vs RM_TempSlot
                ((of_nat CLIENT_SCRATCH_VADDR) + ((of_nat i) << seL4_PageBits))
                seL4_AllRights
   od"
definition move_router_scratch_frame ::
  "(event option) rm_monad"
  where
  "move_router_scratch_frame =
  do i \leftarrow gets (\lambda(tc, rm). var_i rm);
      move_cap RM_Client_CSpace
               ((of_nat CLIENT_CPTR_FIRST_SCRATCH_ITEM) + (of_nat i))
               RM_SelfCSpace RM_TempSlot
   od"
definition set_router_tcb_ip_sp ::
  "(event option) rm_monad"
 where
  "set_router_tcb_ip_sp =
   set_tcb_ip_sp RM_Client_Tcb (of_nat CLIENT_ENTRY_POINT) 0"
definition copy_timer_to_client ::
  "(event option) rm_monad"
 where
  "copy_timer_to_client =
  do ep \leftarrow gets (\lambda(tc, rm). virtualirq_endpoint rm);
```

```
copy_untyped_to_client ::
"(event option) rm_monad"
where
"copy_untyped_to_client =
  do i ← gets (λ(tc, rm). var_i rm);
    item ← allocate_large_item;
    copy_cap RM_Client_CSpace
        ((of_nat CLIENT_CPTR_FIRST_UNTYPED_ITEM) + (of_nat i)) RM_SelfCSpace
        item seL4_AllRights
```

od"

definition

definition

definition

```
revoke_ith_untyped ::
"(event option) rm_monad"
where
"revoke_ith_untyped =
do i ← gets (λ(tc, rm). var_i rm);
free_items' ← gets (λ(tc, rm). free_items rm);
untyped_items' ← gets (λ(tc, rm). untyped_items rm);
revoke_untyped RM_SelfCSpace
((untyped_items'!i)!(unat (free_items'!i)))
od"
```

definition

```
wait_for_message ::
"(event option) rm_monad"
where
"wait_for_message =
  do command_endpoint' ← gets (λ(tc, rm). command_endpoint rm);
    lift_tc_rm (seL4_Wait command_endpoint' 0)
  od"
```

A.6 Router Manager program

```
(********* CASE 0 *********)
              cnd_jmp (\lambda(tc,rm). message rm = 0) 4,
              jmp 13,
(*** TEARDOWN ***)
(*LINE 4*)
              u_instr (modify (\lambda(tc, rm). (tc, rm(|var_i:=0|)))),
              jmp 10,
(*LINE 6*)
              cnd_jmp (\lambda(tc, rm). (free_items rm)!(var_i rm) \geq
                 (total_items rm)!(var_i rm)) 10,
              k_instr revoke_ith_untyped,
              u_instr (modify (\lambda(tc, rm). (tc, free_items_update
                 (λx. x[(var_i rm) := (x ! (var_i rm)) + 1]) rm))),
              jmp 6,
(*LINE 10*)
              cnd_jmp (\lambda(tc, rm). var_i rm < NUM_SIZES) 6,
              k_instr sanitise_network_cards,
              jmp 86,
(********* CASE 1 *********)
(*LINE 13*)
              cnd_jmp (\lambda(tc, rm). message rm = 1) 15,
              jmp 50,
(*** TEARDOWN ***)
(*LINE 15*)
              u_instr (modify (\lambda(tc, rm). (tc, rm(|var_i:=0|)))),
              jmp 21,
(*LINE 17*)
              cnd_jmp (\lambda(tc, rm). (free_items rm)!(var_i rm) \geq
                 (total_items rm)!(var_i rm)) 21,
              k_instr revoke_ith_untyped,
              u_instr (modify (\lambda(tc, rm). (tc, free_items_update
                 (λx. x[(var_i rm) := x!(var_i rm) + 1]) rm))),
              jmp 17,
(*LINE 21*)
              cnd_{jmp} (\lambda(tc, rm). var_i rm < NUM_SIZES) 17,
              k_instr sanitise_network_cards,
(*** SETUP ***)
              k_instr create_router_cspace_object,
              k_instr mutate_router_cspace_cap,
              u_instr (modify (\lambda(tc, rm). (tc, rm(|var_i:=0|)))),
              jmp 31,
(*LINE 27*)
              k_instr create_router_scratch_space,
              k_instr map_router_scratch_frame,
              k_instr move_router_scratch_frame,
              u_instr (modify (\lambda(tc, rm). (tc, rm(|var_i:=var_i rm + 1|)))),
```

(*LINE 31*) cnd_{jmp} ($\lambda(tc, rm)$. var_i rm < CLIENT_NUM_SCRATCH_FRAMES) 27, k_instr create_router_ipc_object, k_instr map_router_ipc_frame, k_instr create_router_tcb_object, k_instr setup_router_tcb, k_instr set_router_tcb_ip_sp, k_instr copy_timer_to_client, u_instr (modify (λ (tc, rm). (tc, rm(|var_i:=0|)))), jmp 42, (*LINE 40*) k_instr copy_untyped_to_client, u_instr (modify (λ (tc, rm). (tc, rm(|var_i:=var_i rm + 1|)))), (*LINE 42*) cnd_jmp (λ (tc, rm). var_i rm < CLIENT_NUM_UNTYPED_ITEMS) 40, k_instr copy_asid_pool_to_client, k_instr copy_nic_A_to_client, k_instr copy_nic_D_to_client, k_instr copy_tcb_cap_to_client, k_instr move_ipc_buffer_cap_to_client, k_instr mutate_cspace_to_client, jmp 86, (********* CASE 2 *********) (*LINE 50*) cnd_jmp (λ (tc, rm). message rm = 2) 52, jmp 86. (***** TEARDOWN *****) (*LINE 52*) u_instr (modify (λ (tc, rm). (tc, rm(|var_i:=0|)))), jmp 58, (*LINE 54*) cnd_jmp (λ (tc, rm). (free_items rm)!(var_i rm) \geq (total_items rm)!(var_i rm)) 58, k_instr revoke_ith_untyped, u_instr (modify (λ (tc, rm). (tc, free_items_update (λx. x[(var_i rm) := x!(var_i rm) + 1]) rm))), jmp 54, (*LINE 58*) cnd_jmp (λ (tc, rm). var_i rm < NUM_SIZES) 54, k_instr sanitise_network_cards,

```
(***** SETUP *****)
              k_instr create_router_cspace_object,
              k_instr mutate_router_cspace_cap,
              u_instr (modify (\lambda(tc, rm). (tc, rm(|var_i:=0|)))),
              jmp 68,
(*LINE 64*)
              k_instr create_router_scratch_space,
              k_instr map_router_scratch_frame,
              k_instr move_router_scratch_frame,
              u_instr (modify (\lambda(tc, rm). (tc, rm(|var_i:=0|)))),
(*LINE 68*)
              cnd_jmp (\lambda(tc, rm). var_i rm < CLIENT_NUM_SCRATCH_FRAMES) 64,
              k_instr create_router_ipc_object,
              k_instr map_router_ipc_frame,
              k_instr create_router_tcb_object,
              k_instr setup_router_tcb,
              k_instr set_router_tcb_ip_sp,
              k_instr copy_timer_to_client,
              u_instr (modify (\lambda(tc, rm). (tc, rm(|var_i:=0|)))),
              jmp 79,
(*LINE 77*)
              k_instr copy_untyped_to_client,
              u_instr (modify (\lambda(tc, rm). (tc, rm(|var_i:=var_i rm + 1|)))),
(*LINE 79*)
              cnd_jmp (\lambda(tc, rm). var_i rm < CLIENT_NUM_UNTYPED_ITEMS) 77,
              k_instr copy_asid_pool_to_client,
              k_instr copy_nic_A_to_client,
              k_instr copy_nic_D_to_client,
              k_instr copy_tcb_cap_to_client,
              k_instr move_ipc_buffer_cap_to_client,
              k_instr mutate_cspace_to_client,
(*LINE 86*)
              jmp 0
  ]"
```

end

Appendix **B**

libseL4 Specification

theory libseL4Interface

imports seL4Interface

begin

B.1 libseL4 Interface

```
definition
```

```
sel4_CNode_Copy ::
"sel4_CPtr ⇒ sel4_CPtr ⇒ sel4_Word ⇒ sel4_CPtr ⇒ sel4_CPtr ⇒ sel4_Word
⇒ sel4_CapRights
⇒ (event option) tc_monad"
where
"sel4_CNode_Copy service dest_index dest_depth src_root src_index src_depth rights=
do sel4_SetMR 0 dest_index;
sel4_SetMR 1 (dest_depth && 0xff);
sel4_SetMR 2 src_index;
sel4_SetMR 3 (src_depth && 0xff);
sel4_SetMR 4 rights;
sel4_SetMR 4 rights;
sel4_Call service (sel4_MessageInfo 5 1 (of_nat (unat CNodeCopy)) 0)
od"
```

```
\label{eq:sel4_CNode_Move ::} $$"sel4_CPtr \Rightarrow sel4_Word \Rightarrow sel4_CPtr \Rightarrow sel4_CPtr \Rightarrow sel4_Word $$
```

```
⇒ (event option) tc_monad"
where
"seL4_CNode_Move service dest_index dest_depth src_root src_index src_depth =
    do seL4_SetCap 0 src_root;
        seL4_SetMR 0 dest_index;
        seL4_SetMR 1 (dest_depth && 0xff);
        seL4_SetMR 2 src_index;
        seL4_SetMR 3 (src_depth && 0xff);
        seL4_Call service (seL4_MessageInfo 4 1 (of_nat (unat CNodeMove)) 0)
        od"
```

```
seL4_Untyped_Retype ::
"seL4_CPtr ⇒ seL4_Word ⇒ seL4_Word ⇒ seL4_CPtr ⇒ seL4_Word ⇒ seL4_Word
⇒ seL4_CPtr ⇒ seL4_Word
⇒ (event option) tc_monad"
where
"seL4_Untyped_Retype service type size_bits root node_index node_depth node_offset
node_window =
do seL4_SetCap 0 root;
seL4_SetMR 1 size_bits;
seL4_SetMR 2 node_index;
seL4_SetMR 3 node_depth;
seL4_SetMR 4 node_offset;
seL4_SetMR 5 node_window;
```

```
seL4_Call service (seL4_MessageInfo 6 1 (of_nat (unat UntypedRetype)) 0)
od"
```

definition

```
seL4_CNode_Revoke ::
"seL4_CNode => seL4_Word => 8 word
=> (event option) tc_monad"
where
"seL4_CNode_Revoke service index depth =
    do seL4_SetMR 0 index;
        seL4_SetMR 1 (of_nat (unat depth) && 0xff);
        seL4_Call service (seL4_MessageInfo 2 0 (of_nat (unat CNodeRevoke)) 0)
    od"
```

```
seL4_CNode_Mutate ::
"seL4_CNode ⇒ seL4_Word ⇒ 8 word ⇒ seL4_CNode ⇒ seL4_Word ⇒ (8 word)
⇒ seL4_CapData
⇒ (event option) tc_monad"
where
"seL4_CNode_Mutate service dest_index dest_depth src_root src_index src_depth
badge =
do seL4_SetCap 0 src_root;
seL4_SetMR 0 dest_index;
seL4_SetMR 1 (of_nat (unat dest_depth) && Oxff);
seL4_SetMR 2 src_index;
seL4_SetMR 3 (of_nat (unat src_depth) && Oxff);
seL4_SetMR 4 badge;
seL4_Call service (seL4_MessageInfo 5 1 (of_nat (unat CNodeMutate)) 0)
od"
```

```
sel4_ARM_Page_Map ::
"sel4_ARM_Page ⇒ sel4_ARM_PageDirectory ⇒ sel4_Word ⇒ sel4_CapRights
⇒ sel4_ARM_VMAttributes
⇒ (event option) tc_monad"
where
"sel4_ARM_Page_Map service pd vaddr rights attr =
    do sel4_SetCap 0 pd;
        sel4_SetMR 0 vaddr;
        sel4_SetMR 1 rights;
        sel4_SetMR 2 attr;
        sel4_Call service (sel4_MessageInfo 3 1 (of_nat (unat ARMPageMap)) 0)
        od"
definition sel4_TCB_Configure ::
"sel4_TCB ⇒ sel4_Word ⇒ 8 word ⇒ sel4_CNode ⇒ sel4_CapData ⇒ sel4_CNode
    ⇒ sel4_CapData ⇒ sel4_Word ⇒ sel4_CPtr
```

```
⇒ (event option) tc_monad"
where
"seL4_TCB_Configure service fault_ep priority cspace_root cspace_root_data
    vspace_root vspace_root_data buffer bufferFrame =
    do seL4_SetCap 0 cspace_root;
        seL4_SetCap 1 vspace_root;
```

```
seL4_SetCap 0 bufferFrame;
```

```
seL4_SetMR 0 fault_ep;
seL4_SetMR 1 (of_nat (unat priority) && Oxff);
seL4_SetMR 2 cspace_root;
seL4_SetMR 3 vspace_root;
seL4_SetMR 4 buffer;
seL4_Call service (seL4_MessageInfo 5 3 (of_nat (unat TCBConfigure)) 0)
od"
```

```
seL4_TCB_WriteRegisters ::
\texttt{"seL4\_TCB} \ \Rightarrow \ \texttt{bool} \ \Rightarrow \ \texttt{8} \ \texttt{word} \ \Rightarrow \ \texttt{seL4\_Word} \ \Rightarrow \ \texttt{seL4\_UserContext}
 \Rightarrow (event option) tc_monad"
where
"seL4_TCB_WriteRegisters service resume_target arch_flags count regs =
   do seL4_SetMR 0
         (of_nat (unat ((of_bl::bool list \Rightarrow 1 word) [resume_target])) && Ox1
         // ((of_nat (unat arch_flags) && Oxff) << 8));</pre>
       seL4_SetMR 1 count;
      seL4_SetMR 2 (upc regs);
      seL4_SetMR 3 (usp regs);
      seL4_SetMR 4 (ucpsr regs);
      seL4_SetMR 5 (ur0 regs);
      seL4_SetMR 6 (ur1 regs);
      seL4_SetMR 7 (ur8 regs);
      seL4_SetMR 8 (ur9 regs);
      seL4_SetMR 9 (ur10 regs);
      seL4_SetMR 10 (ur11 regs);
      seL4_SetMR 11 (ur12 regs);
      seL4_SetMR 12 (ur4 regs);
      seL4_SetMR 13 (ur2 regs);
      seL4_SetMR 14 (ur3 regs);
      seL4_SetMR 15 (ur5 regs);
      seL4_SetMR 16 (ur6 regs);
      seL4_SetMR 17 (ur7 regs);
      seL4_SetMR 18 (ur14 regs);
      seL4_Call service (seL4_MessageInfo 19 0 (of_nat (unat TCBWriteRegisters)) 0)
    od"
```

end

```
theory seL4Interface
```

imports libseL4Functions

begin

B.2 seL4 Interface

```
definition seL4_Wait ::
  "seL4_CPtr \Rightarrow seL4_Word \Rightarrow (event option) tc_monad"
  where
  "seL4_Wait src sender =
     do modify (\lambda(i,u). (i, u (RO := src)));
         modify (\lambda(i,u). (i, u (R7 := seL4_SysWait)));
         return (Some (SyscallEvent SysWait))
     od"
definition seL4_Call ::
  "seL4_CPtr \Rightarrow seL4_MessageInfo \Rightarrow (event option) tc_monad"
  where
  "seL4_Call dest msgInfo =
     do modify (\lambda(i,u). (i, u (RO := dest)));
         modify (\lambda(i,u). (i, u (R1 := msgInfo)));
         MRO \leftarrow gets (seL4_GetMR O);
         modify (\lambda(i,u). (i, u (R2 := MR0)));
         MR1 \leftarrow gets (seL4\_GetMR 1);
         modify (\lambda(i, u). (i, u (R3 := MR1)));
         MR2 \leftarrow gets (seL4\_GetMR 2);
         modify (\lambda(i,u). (i, u (R4 := MR2)));
         MR3 \leftarrow gets (seL4_GetMR 3);
         modify (\lambda(i,u). (i, u (R5 := MR3)));
         modify (\lambda(i,u). (i, u (R7 := seL4_SysCall)));
         return (Some (SyscallEvent SysCall))
     od"
```

end theory libseL4Functions imports libseL4Types begin

B.3 libseL4 Functions

```
definition seL4_GetMR ::
  "nat \Rightarrow thread_abstract_context \Rightarrow seL4_Word"
  where
  "seL4_GetMR n = (\lambda(i, u). (ipc_msg i) ! n)"
definition seL4_SetMR ::
  "nat \Rightarrow seL4_Word \Rightarrow (event option) tc_monad"
  where
  "seL4_SetMR n w =
     do modify (\lambda(i, u). ((ipc_msg_update (\lambdax. x[n := w]) i, u)));
         return None
     od"
definition seL4_SetCap ::
  "nat \Rightarrow seL4_CPtr \Rightarrow (event option) tc_monad"
  where
  "seL4_SetCap n w =
     do modify (\lambda(i, u). ((ipc_caps_update (\lambda x. x[n := w]) i, u)));
         return None
      od"
end
theory libseL4Types
imports Main Event_H UserADTs
begin
        libseL4 Types
B.4
types seL4_Word = word32
types seL4_CPtr = seL4_Word
```

```
definition seL4_PageBits :: "nat" where "seL4_PageBits = 12"
definition seL4_SlotBits :: "nat" where "seL4_SlotBits = 4"
definition seL4_TCBBits :: "nat" where "seL4_TCBBits = 9"
```

definition seL4_WordBits :: "nat" where "seL4_WordBits = 32"

```
definition seL4_EndpointBits :: "nat" where "seL4_EndpointBits = 4"
definition seL4_PageTableBits :: "nat" where "seL4_PageTableBits = 10"
definition seL4_PageDirBits :: "nat" where "seL4_PageDirBits = 14"
definition seL4_ASIDPoolBits :: "nat" where "seL4_ASIDPoolBits = 12"
```

```
definition seL4_Frame_Args :: "nat" where "seL4_Frame_Args = 4"
definition seL4_Frame_MRs :: "nat" where "seL4_Frame_MRs = 7"
definition seL4_Frame_HasNPC :: "nat" where "seL4_Frame_HasNPC = 0"
```

```
types seL4_ARM_Page = seL4_CPtr
types seL4_ARM_PageTable = seL4_CPtr
types seL4_ARM_PageDirectory = seL4_CPtr
types seL4_ARM_ASIDControl = seL4_CPtr
types seL4_ARM_ASIDPool = seL4_CPtr
```

```
record seL4_UserContext =
```

upc :: seL4_Word usp :: seL4_Word ucpsr :: seL4_Word ur0 :: seL4_Word ur1 :: seL4_Word ur8 :: seL4_Word ur9 :: seL4_Word ur10 :: seL4_Word ur11 :: seL4_Word ur12 :: seL4_Word ur4 :: seL4_Word ur2 :: seL4_Word ur3 :: seL4_Word ur5 :: seL4_Word ur6 :: seL4_Word ur7 :: seL4_Word

ur14 :: seL4_Word

```
types seL4_ARM_VMAttributes = seL4_Word
definition seL4_ARM_PageCacheable :: "seL4_Word" where "seL4_ARM_PageCacheable = 1"
```

```
definition seL4_ARM_ParityEnabled :: "seL4_Word" where "seL4_ARM_ParityEnabled = 2"
definition seL4_ARM_Default_VMAttributes :: "seL4_Word" where
    "seL4_ARM_Default_VMAttributes = 3"
```

```
types seL4_Syscall_ID = seL4_Word
definition seL4_SysSend :: "seL4_Word" where "seL4_SysSend = -1"
definition seL4_SysNBSend :: "seL4_Word" where "seL4_SysNBSend = -2"
definition seL4_SysCall :: "seL4_Word" where "seL4_SysCall = -3"
definition seL4_SysWait :: "seL4_Word" where "seL4_SysWait = -4"
definition seL4_SysReply :: "seL4_Word" where "seL4_SysReply = -5"
definition seL4_SysReplyWait :: "seL4_Word" where "seL4_SysReply = -5"
definition seL4_SysReplyWait :: "seL4_Word" where "seL4_SysReplyWait = -6"
definition seL4_SysVield :: "seL4_Word" where "seL4_SysVield = -7"
definition seL4_SysDebugPutChar :: "seL4_Word" where "seL4_SysDebugPutChar = -8"
definition seL4_SysDebugHalt :: "seL4_Word" where "seL4_SysDebugHalt = -9"
definition seL4_SysDebugSnapshot :: "seL4_Word" where "seL4_SysDebugSnapshot = -10"
definition seL4_SysDebugRun :: "seL4_Word" where "seL4_SysDebugRun = -11"
```

types seL4_CapData = seL4_Word

definition

```
seL4_CapData ::
"5 word ⇒ 18 word ⇒ seL4_CapData"
where
"seL4_CapData guard_size guard_bits =
      ((of_nat (unat (0::3 word))) << 0)
  + ((of_nat (unat guard_size)) << 3)
  + ((of_nat (unat guard_bits)) << 8)"</pre>
```

types seL4_MessageInfo = seL4_Word

definition

seL4_MessageInfo :: "7 word ⇒ 2 word ⇒ 3 word ⇒ 20 word ⇒ seL4_MessageInfo" where "seL4_MessageInfo 1 extraCaps capsUnwrapped label = ((of_nat (unat 1)) << 0)</pre>

```
+ ((of_nat (unat extraCaps)) << 7)
+ ((of_nat (unat capsUnwrapped)) << 9)
+ ((of_nat (unat label)) << 12)"</pre>
```

```
types seL4_ObjectType = seL4_Word
definition seL4_UntypedObject :: "seL4_Word" where "seL4_UntypedObject = 0"
definition seL4_TCBObject :: "seL4_Word" where "seL4_TCBObject = 1"
definition seL4_EndpointObject :: "seL4_Word" where "seL4_EndpointObject = 2"
definition seL4_AsyncEndpointObject ::
    "seL4_Word" where "seL4_AsyncEndpointObject = 3"
definition seL4_CapTableObject :: "seL4_Word" where "seL4_CapTableObject = 4"
definition seL4_NonArchObjectTypeCount :: "seL4_Word" where
    "seL4_NonArchObjectTypeCount = 5"
```

```
types seL4_CapRights = seL4_Word
definition seL4_CanWrite :: "seL4_Word" where "seL4_CanWrite = 1"
definition seL4_CanRead :: "seL4_Word" where "seL4_CanRead = 2"
definition seL4_CanGrant :: "seL4_Word" where "seL4_CanGrant = 4"
definition seL4_AllRights :: "seL4_Word" where "seL4_AllRights = 7"
definition seL4_Transfer_Mint :: "seL4_Word" where "seL4_Transfer_Mint = 256"
```

```
datatype seL4_IPCBuffer =
    seL4_IPCBuffer seL4_MessageInfo "(seL4_Word list)" seL4_Word
    "(seL4_CPtr list)" seL4_CPtr seL4_CPtr seL4_Word
```

```
\begin{array}{ll} \texttt{ipc\_msg} :: \\ \texttt{"seL4\_IPCBuffer} \Rightarrow \texttt{seL4\_Word} \texttt{list"} \\ \texttt{where} \\ \texttt{"ipc\_msg} = \\ & (\lambda \texttt{x}. \ (\texttt{case x of} \\ & \quad \texttt{seL4\_IPCBuffer tag msg userData caps receiveCNode} \\ & \quad \texttt{receiveIndex receiveDepth} \\ & \quad \Rightarrow \texttt{msg}))\texttt{"} \end{array}
```

```
\begin{split} & \text{ipc\_msg\_update ::} \\ & \text{"(seL4\_Word list} \Rightarrow \text{seL4\_Word list)} \Rightarrow \text{seL4\_IPCBuffer} \Rightarrow \text{seL4\_IPCBuffer"} \\ & \text{where} \\ & \text{"ipc\_msg\_update f =} \\ & (\lambda x. \text{ (case x of} \\ & \text{seL4\_IPCBuffer tag msg userData caps receiveCNode} \\ & \text{ receiveIndex receiveDepth} \\ & \Rightarrow \text{seL4\_IPCBuffer tag (f msg) userData caps receiveCNode} \\ & \text{ receiveIndex receiveDepth}) \end{split}
```

definition

```
ipc\_caps ::
"seL4_IPCBuffer \Rightarrow seL4_Word list"
where
"ipc_caps =
(\lambda x. (case x of
seL4_IPCBuffer tag msg userData caps receiveCNode
receiveIndex receiveDepth
\Rightarrow caps))"
```

definition

```
\begin{split} & \text{ipc_caps_update ::} \\ & \text{"(seL4_Word list} \Rightarrow \text{seL4_Word list)} \\ & \Rightarrow \text{seL4_IPCBuffer} \\ & \Rightarrow \text{seL4_IPCBuffer"} \\ & \text{where} \\ & \text{"ipc_caps_update f =} \\ & (\lambda x. \text{ (case x of} \\ & \text{ seL4_IPCBuffer tag msg userData caps receiveCNode} \\ & \text{ receiveIndex receiveDepth} \\ & \Rightarrow \text{ seL4_IPCBuffer tag msg userData (f caps) receiveCNode} \\ & \text{ receiveIndex receiveDepth}))" \end{split}
```

types seL4_CNode = seL4_CPtr
types seL4_IRQHandler = seL4_CPtr
types seL4_IRQControl = seL4_CPtr
types seL4_TCB = seL4_CPtr
types seL4_Untyped = seL4_CPtr

```
definition InvalidInvocation :: "seL4_Word" where "InvalidInvocation = 0"
definition UntypedRetype :: "seL4_Word" where "UntypedRetype = 1"
definition TCBReadRegisters :: "seL4_Word" where "TCBReadRegisters = 2"
definition TCBWriteRegisters :: "seL4_Word" where "TCBWriteRegisters = 3"
definition TCBCopyRegisters :: "seL4_Word" where "TCBCopyRegisters = 4"
definition TCBConfigure :: "seL4_Word" where "TCBConfigure = 5"
definition TCBSetPriority :: "seL4_Word" where "TCBSetPriority = 6"
definition TCBSetIPCBuffer :: "seL4_Word" where "TCBSetIPCBuffer = 7"
definition TCBSetSpace :: "seL4_Word" where "TCBSetSpace = 8"
definition TCBSuspend :: "seL4_Word" where "TCBSuspend = 9"
definition TCBResume :: "seL4_Word" where "TCBResume = 10"
definition CNodeRevoke :: "seL4_Word" where "CNodeRevoke = 11"
definition CNodeDelete :: "seL4_Word" where "CNodeDelete = 12"
definition CNodeRecycle :: "seL4_Word" where "CNodeRecycle = 13"
definition CNodeCopy :: "seL4_Word" where "CNodeCopy = 14"
definition CNodeMint :: "seL4_Word" where "CNodeMint = 15"
definition CNodeMove :: "seL4_Word" where "CNodeMove = 16"
definition CNodeMutate :: "seL4_Word" where "CNodeMutate = 17"
definition CNodeRotate :: "seL4_Word" where "CNodeRotate = 18"
definition CNodeSaveCaller :: "seL4_Word" where "CNodeSaveCaller = 19"
definition IRQIssueIRQHandler :: "seL4_Word" where "IRQIssueIRQHandler = 20"
definition IRQInterruptControl :: "seL4_Word" where "IRQInterruptControl = 21"
definition IRQAckIRQ :: "seL4_Word" where "IRQAckIRQ = 22"
definition IRQSetIRQHandler :: "seL4_Word" where "IRQSetIRQHandler = 23"
definition IRQClearIRQHandler :: "seL4_Word" where "IRQClearIRQHandler = 24"
definition nInvocationLabels :: "seL4_Word" where "nInvocationLabels = 25"
```

```
definition ARMPageTableMap :: "seL4_Word" where "ARMPageTableMap = nInvocationLabels"
definition ARMPageTableUnmap :: "seL4_Word" where
    "ARMPageTableUnmap = nInvocationLabels + 1"
definition ARMPageMap :: "seL4_Word" where "ARMPageMap = nInvocationLabels + 2"
definition ARMPageRemap :: "seL4_Word" where "ARMPageRemap = nInvocationLabels + 3"
definition ARMPageUnmap :: "seL4_Word" where "ARMPageUnmap = nInvocationLabels + 4"
definition ARMPageFlushCaches :: "seL4_Word" where
    "ARMPageFlushCaches = nInvocationLabels + 5"
```

```
definition ARMASIDControlMakePool :: "seL4_Word" where
  "ARMASIDControlMakePool = nInvocationLabels + 6"
definition ARMASIDPoolAssign :: "seL4_Word" where
  "ARMASIDPoolAssign = nInvocationLabels + 7"
```

types thread_abstract_context = "seL4_IPCBuffer × user_context"
types 'a tc_monad = "(thread_abstract_context, 'a) nondet_monad"

end

Appendix C

System Process

theory Decode_I

imports RouterManager_I

begin

C.1 Router Manager Functional specification system process

C.1.1 Lifting from rm_monad to user_state transitions

```
consts
  decode_sac_state_I ::
   "user_mem ⇒ sac_abstract_state"
consts
  decode_ipc_I ::
   "user_mem ⇒ seL4_IPCBuffer"
definition
  decode_I ::
   "user_state ⇒ thread_abstract_context × sac_abstract_state"
  where
   "decode_I = (λ(uc, ip, um).
      ((decode_ipc_I um, uc), ip_update (λx. ip >> 2) (decode_sac_state_I um))) "
```

consts

```
\label{eq:encode_sac_state_I} \encode\_sac\_state\_I :: $$"seL4_IPCBuffer \Rightarrow sac_abstract\_state \Rightarrow user_mem"$$
```

```
encode_I ::

"thread_abstract_context \times sac_abstract_state \Rightarrow user_state"

where

"encode_I = (\lambda((ipc, uc), s). (uc, (ip s << 2), encode_sac_state_I ipc s))"
```

definition

```
encode_return_I ::

"((event option) × thread_abstract_context × sac_abstract_state) set × bool

\Rightarrow (user_state × event option) set"

where

"encode_return_I = (\lambda(Q,b). {(encode_I s, e)/s e. (e, s) \in Q})"
```

definition

```
\begin{split} & \text{lift_I}:: \\ & \text{"instr} \Rightarrow (\texttt{user\_state} \Rightarrow (\texttt{user\_state} \times \texttt{event option}) \texttt{ set}) \texttt{"} \\ & \text{where} \\ & \text{"lift_I} \texttt{ i = } (\lambda\texttt{s. encode\_return\_I} \texttt{ (i (decode\_I \texttt{s})))"} \end{split}
```

C.2 User state translations

consts program_text_I :: "word32 ⇒ word8 option"

consts

```
project_program_text ::
"user_state \Rightarrow (word32 \rightarrow word8)"
```

consts

consts

axioms decode_router_manager_program_text:

"decode_program_text program_text_I = map lift_I router_manager"

C.3 The user specification

definition

```
uop_I ::
"user_state ⇒ (user_state × event option) set"
where
"uop_I us =
  ((decode_program_text (project_program_text us))!(project_program_counter us)) us"
```

C.4 The system process

definition

```
ADT_I :: "(state global_state, observable, global_transition) data_type"
where
"ADT_I ≡
(| Init = λs. Init_A, Fin = Fin_A,
Step = global_automaton (do_user_op uop_I) kernel_call_A|)"
```

end

theory UserADTs

imports Simulation Invariant

begin

C.5 System process

We constrain idle thread behaviour, so we distinguish three main machine control states:

```
datatype mode = UserMode | KernelMode | IdleMode
```

```
types user_mem = "word32 ⇒ word8 option"
types user_state = "user_context × word32 × user_mem"
```

```
types observable_user_state = "user_context × user_mem"
types machine_no_mem = "irq set × machine_state_rest"
```

We take as observable state the user state (register context and memory), the internal machine state without memory, the current-thread pointer in both specs (i.e. *cur_thread* in the abstract kernel model and *ksCurThread* in the executable specification), the mode, and the current event (if any).

Currently, user memory is all memory accessible to any user (everything of type UserData). A virtual-memory view could be built on top of this.

```
types observable = "observable_user_state \times machine_no_mem \times word32 \times mode \times event option"
```

The global state contains the current register context of the machine as well as the internal kernel state, the mode and the current event (if any).

```
types 'k global_state = "(user_context \times word32 \times 'k) \times mode \times event option"
```

The global-transition type reflects the operations of the abstract data types.

datatype

```
global_transition =
    KernelTransition
    UserTransition
    UserEventTransition
    IdleEventTransition
```

C.6 The global_automaton

The following definition models machine, user, and kernel entry/exit behaviour abstractly. It also constrains the behaviour of the idle thread. The first parameter specifies all possible user transition on the kernel state (which is expected to inject the new user memory into the kernel state) and the second parameter provides a kernel specification (it will later be used with the abstract as well as the executable specification).

The first and second transition sets are kernel calls being serviced by the kernel. In the first set the kernel exits into UserMode, in the second the kernel exits into IdleMode.

The third transition set is normal user execution.

The fourth transition set is user execution that generates a syscall event and traps into the kernel.

The fifth transition set is user execution being interrupted by an interrupt event.

The sixth transition set is idle execution.

```
global_automaton ::
  "((user_context 	imes word32 	imes 'k) 	imes (user_context 	imes word32 	imes 'k)
    \times mode \times event option) set
   \Rightarrow (event \Rightarrow ((user_context \times 'k) \times mode \times (user_context \times 'k)) set)
   \Rightarrow global_transition \Rightarrow ('k global_state 	imes 'k global_state) set"
where
"global_automaton user_transition kernel_call t \equiv case t of
  KernelTransition \Rightarrow
     { ( ((uc, ip, s), KernelMode, Some e),
          ((uc', uc' LR_svc, s'), UserMode, None) ) |uc ip s uc' s' e.
            ((uc, s), UserMode, (uc', s')) \in kernel_call e \}
   \cup { ( ((uc, ip, s), KernelMode, Some e),
          ((uc', ip', s'), IdleMode, None) ) |uc ip s uc' ip' s' e.
            ((uc, s), IdleMode, (uc', s')) \in kernel_call e }
| UserTransition \Rightarrow
     { ( (g, UserMode, None), (g', UserMode, e) ) |g g' e.
     (g, g', UserMode, e) ∈ user_transition}
   \cup { ( ((uc, ip, g), UserMode, None),
          ((uc'(FaultInstruction := ip'), ip', g'), KernelMode, Some ev) )
          /uc ip g uc' ip' g' ev.
            ((uc, ip, g), (uc', ip', g'), KernelMode, Some ev) \in user_transition}
   \cup { ( (g, UserMode, None), (g, KernelMode, Some Interrupt) ) |g. True }
| IdleEventTransition \Rightarrow
   { ( (g, IdleMode, None), (g, KernelMode, Some Interrupt) ) |g. True }"
```

After kernel initialisation, the machine is in UserMode, running the initial thread.

definition

```
Init_A :: "state global_state set"
where
   "Init_A = {((empty_context, 0, init_A_st), UserMode, None)}"
```

A pointer is inside a user frame if its top bits point to any object of type IntData

definition

"in_user_frame p s $\equiv \exists$ sz. typ_at (AArch (AIntData sz)) (p && ~~ mask (pageBitsForSize sz)) s"

The content of user memory is stored in the machine state. Only locations tagged as inside a user frame are accessible.

definition

"user_mem s $\equiv \lambda p$.

```
if in_user_frame p s
then Some (underlying_memory (machine_state s) p)
else None"
```

```
"user_memory_update um \equiv modify (\lambdams.
ms(|underlying_memory := (\lambdaa. case um a of Some x \Rightarrow x
/ None \Rightarrow underlying_memory ms a)|))"
```

definition

```
do_user_op' ::

" (user_state ⇒ (user_state × event option) set)

⇒ user_context × word32

⇒ ((user_context × word32) × mode × event option) s_monad"

where

"do_user_op' uop ≡ (λ(uc, ip).

do um ← gets user_mem;

    ((uc',ip',um'), e) ← select (uop (uc,ip,um));

    m ← (case e of None ⇒ return UserMode | Some ev ⇒ return KernelMode);

    do_machine_op (user_memory_update um');

    return ((uc',ip'), m, e)

    od)"
```

definition

```
do_user_op ::
" (user_state ⇒ (user_state × event option) set)
⇒ ((user_context × word32 × state) × ((user_context × word32 × state) × mode
            × event option)) set"
where
"do_user_op uop ≡
        {((uc, ip, s), ((uc', ip', s'), m, e)).
            (((uc', ip'), m, e), s') ∈ fst (split (do_user_op' uop) ((uc, ip), s)) }"
```

Kernel calls are described completely by the abstract and concrete spec. We model kernel entry by allowing an arbitrary user (register) context. The mode after a kernel call is either user or idle (see also thm in Refine.thy).

```
kernel_entry :: "event \Rightarrow user_context \Rightarrow user_context s_monad" where 
 "kernel_entry e tc \equiv do
```

```
t ← gets cur_thread;
thread_set (λtcb. tcb (| tcb_context := tc |)) t;
call_kernel e;
t' ← gets cur_thread;
thread_get tcb_context t'
od"
```

Extracting the observable state:

definition

"no_mem_machine m \equiv (irq_masks m, machine_state_rest m)"

The final observable state:

definition Fin_A :: "state global_state \Rightarrow observable" where "Fin_A $\equiv \lambda((uc, ip, s), m, e)$. ((uc, user_mem s), no_mem_machine (machine_state s), cur_thread s, m, e)"

Lifting a state relation on kernel states to global states.

definition

"lift_state_relation sr \equiv { (((tc,s),m,e), ((tc,s'),m,e))/s s' m e tc. (s,s') \in sr }"

lemma lift_state_relationD:

"(((tc, s), m, e), ((tc', s'), m', e')) \in lift_state_relation R \implies (s,s') \in R \land tc' = tc \land m' = m \land e' = e" by (simp add: lift_state_relation_def)

lemma lift_state_relationI:

"(s,s') $\in \mathbb{R} \implies$ (((tc, s), m, e), ((tc, s'), m, e)) \in lift_state_relation \mathbb{R} " by (fastsimp simp: lift_state_relation_def)

```
lemma in_lift_state_relation_eq:
  "(((tc, s), m, e), (tc', s'), m', e') ∈ lift_state_relation R ↔
  (s, s') ∈ R ∧ tc' = tc ∧ m' = m ∧ e' = e"
by (auto simp add: lift_state_relation_def)
```

end

Bibliography

- J. Andronick, D. Greenaway, and K. Elphinstone. Towards proving security in the presence of large untrusted components. In G. Klein, R. Huuck, and B. Schlich, editors, *Proceedings of the* 5th Workshop on Systems Software Verification, Vancouver, Canada, Oct 2010. USENIX.
- [2] D. Elkaduwe, G. Klein, and K. Elphinstone. Verified protection model of the seL4 microkernel. In J. Woodcock and N. Shankar, editors, *Proceedings of Verified Software: Theories, Tools and Experiments 2008*, volume 5295 of *Lecture Notes in Computer Science*, pages 99–114, Toronto, Canada, Oct 2008. Springer-Verlag.
- [3] G. Klein, J. Andronick, K. Elphinstone, G. Heiser, D. Cock, P. Derrin, D. Elkaduwe, K. Engelhardt, R. Kolanski, M. Norrish, T. Sewell, H. Tuch, and S. Winwood. seL4: Formal verification of an OS kernel. *Communications of the ACM*, 53(6):107–115, Jun 2010.
- [4] G. Klein, T. Murray, P. Gammie, T. Sewell, and S. Winwood. Provable security: How feasible is it? In *Proceedings of the 13th Workshop on Hot Topics in Operating Systems*, Napa, CA, USA, May 2011.
- [5] G. Klein, T. Sewell, and S. Winwood. *Refinement in the formal verification of seL4*, pages 323–339. Springer, 1st edition, March 2010.
- [6] L. C. Paulson. The isabelle reference manual, 2008.
- [7] S. Winwood, G. Klein, T. Sewell, J. Andronick, D. Cock, and M. Norrish. Mind the gap: A verification framework for low-level C. In S. Berghofer, T. Nipkow, C. Urban, and M. Wenzel, editors, *Proceedings of the 22nd International Conference on Theorem Proving in Higher Order Logics*, volume 5674 of *Lecture Notes in Computer Science*, pages 500–515, Munich, Germany, Aug 2009. Springer-Verlag.