

Securing a Multiprocessor KVM Hypervisor with 🦣 Rust (socc '24)

Yu-Hsun Chiang, Wei-Lin Chang, Shih-Wei Li, Jan-Ting Tu @ National Taiwan University

paper slides created by Yulong Han @ Peking University



Background: Rust programming language



- low-level programming with **no GC**(Garbage Collection)
- emphasizing *performance*, *type-safety* and *concurrency*
- ownership model:
 - one value can only have one owner
 - reference can borrow the value as mutable/immutable
 - immutable references can have many, but only one mutable reference should exist
- variable lifetime check and notation
- fearless **concurrency** in [std] environment



- 1. fn main() {
- 2. let s1 = String::from("Hello, Rust!");
- 3. let s2 = s1;
- 4. // println!("{}", s1); // ERROR: s1 is no longer valid

variable is immutable by default



```
1. fn main() {
2. let s1 = String::from("Hello, Rust!");
3. let s2 = s1;
4. // println!("{}", s1); // ERROR: s1 is no longer valid
5.
6. let s3 = &s2; // Borrow s2
7. println!("Immutable reference: {}", s3);
```

one variable can have multiple immutable references 📃







```
1. fn main() {
        let s1 = String::from("Hello, Rust!");
 2.
        let s_{2} = s_{1};
 3.
        // println!("{}", s1); // ERROR: s1 is no longer valid
 4.
 5.
 6.
        let s3 = &s2; // Borrow s2
 7.
        println!("Immutable reference: {}", s3);
 8.
        let mut s4 = String::from("Mutable Rust");
 9.
10.
        let s5 = \&mut s4;
        s5.push str(" is cool!");
11.
12.
        println!("Mutable reference: {}", s5);
```

only one mutable referece should exist (even s4 is invalid after s5)





```
1. fn main() {
        let s1 = String::from("Hello, Rust!");
 2.
       let s_{2} = s_{1};
 3.
        // println!("{}", s1); // ERROR: s1 is no longer valid
 4.
 5.
 6.
        let s3 = &s2; // Borrow s2
 7.
        println!("Immutable reference: {}", s3);
 8.
 9.
        let mut s4 = String::from("Mutable Rust");
10.
        let s5 = \&mut s4;
11.
        s5.push str(" is cool!");
12.
        println!("Mutable reference: {}", s5);
13.
                                                 functions take and return ownership
        14.
15.
        println!("After function: {}", s6);
16. }
17.
18. fn takes and gives back(s: String) -> String {
                                                                                      Standard Output
        println!("Taking ownership of: {}", s);
19.
                                                            Immutable reference: Hello, Rust!
        s 💽
20.
                                                            Mutable reference: Mutable Rust is cool!
                                                            Taking ownership of: Hello, Rust!
21. }
                                                            After function: Hello, Rust!
```

one more thing: Rust traits





one more thing: Rust traits

```
1. trait Greet {
 2. fn greet(&self) -> String;
 3. }
 4.
 5. struct Person {
      name: String,
 6.
 7. }
 8.
 9. impl Greet for Person {
10. fn greet(&self) -> String {
11. format!("Hello, {}!", self.name)
12. }
13. }
14.
15. fn main() {
                                    Alice
17. name: String::from("Alice"), IT
18. };
     println!("{}", person.greet());
                                         Hello, Alice!
19.
20. }
```





Background: System Virtualization and Hypervisors

- System virtualization: provides multiple virtual hardware on a single hardware
- Key implementations:
 - **CPU** virtualization: vCPU
 - Memory virtualization: memory isolation between virtual machines(VMs)
 - **I/O** virtualization: pass-through, virtio, ...
- Hypervisors(Virtual Machine Monitors, VMMs):
 - the software that provides the virtual HAL for VMs
 - Type-1 hypervisors (bare-metal)
 - Xen, Hyper-V
 - Type-2 hypervisors (as OS applications)
 - VMware Workstation, Virtualbox



UNIVER SUNIVERSITY

Background: KVM (Kernel-based Virtual Machine)

- Part of the Linux kernel (built-in or as kernel modules)
- Implemented by each architecture's code(x86, arm64, riscv, ...)
 - KVM uses ISA's virtualization extension(VT-x, ARM's virtualization extension, RISCV's H extension)
 - KVM uses linux's scheduler to run VMs' vCPUs as Linux processes
- ARM64's privilege levels:
 - ELO: apps
 - EL1 : os or VM os
 - EL2 : hypervisors
 - EL3 : secure monitor
- Highvisor/Lowvisor on ARM64's KVM:
 - highvisor runs in EL1 and take advantage of Linux's kernel
 - lowvisor runs in EL2 as a hypervisor
- in userspace, **QEMU** is used to interact with KVM







Challenge 1: problems of Mutex in Rust

- std::sync::Mutex(std), spin::Mutex(no_std)
- self-deadlocks
 - lock(A)
 - ... // maybe later or in another function
 lock(A) // waits forever!
- lock ordering problems
 - CPU0: lock(A) ... lock(B) // wait for CPU1
 - CPU1: lock(B) ... lock(A) // wait for CPU0



Challenge 2: raw pointers in unsafe Rust

- unsafe is common in low level Rust system programming:
 - pagetable walking
 - memory allocation
 - atomic instructions/inline asm
 - memory-mapped I/O
 - • •
- can we wrap the **raw pointers** to safely access physical memory regions?



Challenge 3: big TCB(Trusted Computing Base) of KVM

- attackers can exploit bugs in KVM and Linux kernel to steal precious data from VMs.
- since KVM uses Linux kernel functionality like scheduling vCPUs, the TCB is the whole kernel.





KrustVM: KVM, but with a new Rust code core

Contributions

- designed KMutex to avoid self-deadlocks and ensure lock ordering
- designed Safe Pointers to guard raw pointer accesses
- rewrite KVM as KrustVM with some functionality in Rust as TCB (named Rcore)



KMutex



avoid self-deadlocks

```
impl<T: ?Sized> KMutex<T> {
    pub fn lock(&mut self)->Guard<'_,Self>{...}
}
```

- &mut will force Rust compiler to check if the Guard has not been dropped, no one can access KMutex
- this is because in Rust there can only exist **one** mutable reference in the variable's lifetime.

KMutex(cont')

- force lock ordering
- manually designed a lock order graph for all Rcore types
- use trait CanGet* for ordering, for example:
 - LEntry, VM info, PMEM info will be implemented CanGetS2PTInfo trait
- the graph has no loop, so all locks will be forced in order



Figure: the order of getting locks

- 1. https://quinedot.github.io/rust-learning/st-reborrow.html
- 2. https://github.com/rust-lang/reference/issues/788

```
1 pub unsafe trait CanGetA {}
2 // SAFETY: We've manually verified the order
3 unsafe impl CanGetA for B {}
5 pub fn get a<T:CanGetA>(:&mut T)->&mut KMutex<A>{...}
6
7 fn foo(ref b: &mut KMutex<B>) {
    let mut b = ref_b.lock(); b: &mut B
9
10
    /* b in the following line gets converted to
       "&mut *b" by the compiler
11
    let ref_a = get_a(b);
12
                            reborrowed here [line 12] as anonymous &mut
13
                            and should live as long as ref_a [line 12-19]
14
    /* this does not compi
       because ref_a's lifetime is not over */
15
16
    let ref_c = get_c(b);
17
    let mut a = ref a.lock();
18
19
    a.do_a_work();
20
    // we can use b after this B requests lock first,
21
22
    b.do_b_work();
                                 but can only be used after anonymous &mut
23 }
                                 of B(by ref a) is dead \odot
```



SUNIVER MARKET I 8 9 B

Safe Pointers

- partition physical memory as regions
- manually define *Pointers* for these regions with range check



KVM host is untrustable

- Rcore unmap itself, all VM's S2PT and host's S2PT from host
 - S2PT(Stage-2 Pagetables) translate VM's guest Physical Address to real Physical Address

DDR

SRAM ROM

Stage 2

- used for memory isolation
- we don't trust KVM's host so we force it uses trap-in to prevent direct physical memory access (making it avisite av
- Therefore, KVM interacts with Reore with Hypercalls
 - hypercalls are requests from Virtual machines' OS to hypervisor
 - just like syscalls from userspace to OS

Application

Rcore supports VM boot image verification by Ed25519 algorithm

Evalution

- benchmarks: hackbench, netperf, ApacheBench, memcached, redis
- Comparable to KVM with max 10% overhead

Figure 4: Application Benchmark Performance.

My comments on the paper

- reminder of previous paper presentation: DRust: Language-Guided Distributed Shared Memory with Fine Granularity, Full Transparency, and Ultra Efficiency
- KrustVM and DRust both take advantage of **Rust**:
 - Leveraging Rust's ownership and lifetime check for (distributed) system programming
 - and more?
- KrustVM's Rcore uses a **formal verification** tool to verify part of its code:
 - Formal verification has become more and more important for system software

Thanks for Watching

paper slides created by Yulong Han @ Peking University