

# CONVEROS: Practical Model Checking for Verifying Rust OS Kernel Concurrency

Ruize Tang<sup>1,\*</sup>, Minghua Wang<sup>2</sup>, Xudong Sun<sup>3</sup>, Lin Huang<sup>2</sup>, Yu Huang<sup>1</sup>, Xiaoxing Ma<sup>1</sup>

<sup>1</sup> Nanjing University

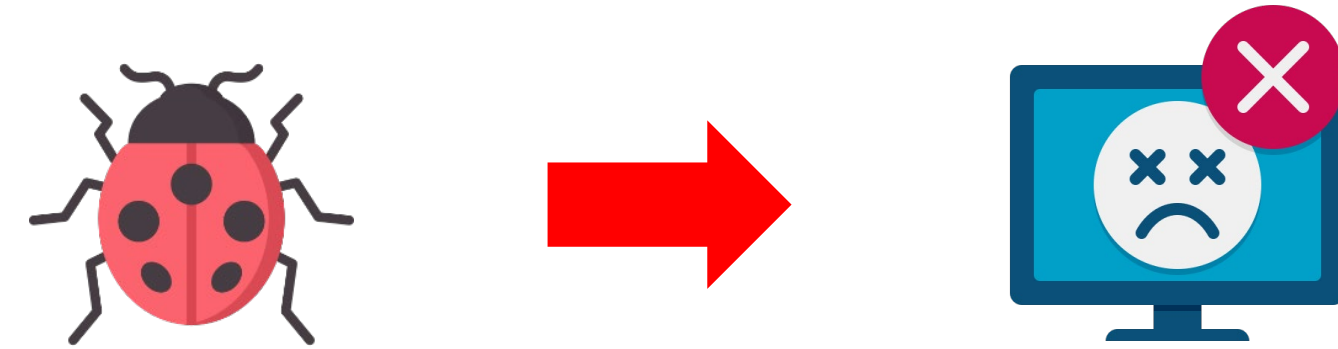
<sup>2</sup> Ant Group

<sup>3</sup> University of Illinois Urbana-Champaign

\* Work done during internship at Ant Group



# Background and motivation



- **Asterinas OS kernel**

- Rust-based, open-source kernel
- Linux ABI compatibility
- 100K LoC, actively developed
- Designed for high reliability

- **Problem: concurrency bugs**

- Persist even in Rust (e.g., deadlocks and data races from design flaws)
- Lead to kernel panics, hangs, etc.

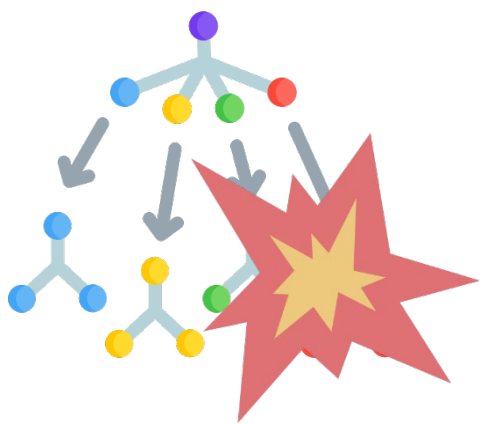
- **Approach: model checking**

- Explores all states of concurrent systems
- Proven effective in detecting deep bugs (e.g., AWS, Microsoft, MongoDB)

# Challenges in applying model checking to OS kernels

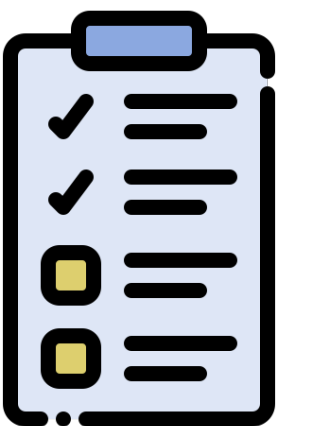
## Challenge 1: Writing high-quality specs

- **High barrier to entry:** Concurrency modules are tricky, requiring both domain knowledge and formal modeling expertise
- **Scalability issues:** Poorly structured specs cause state explosion and limit verification depth



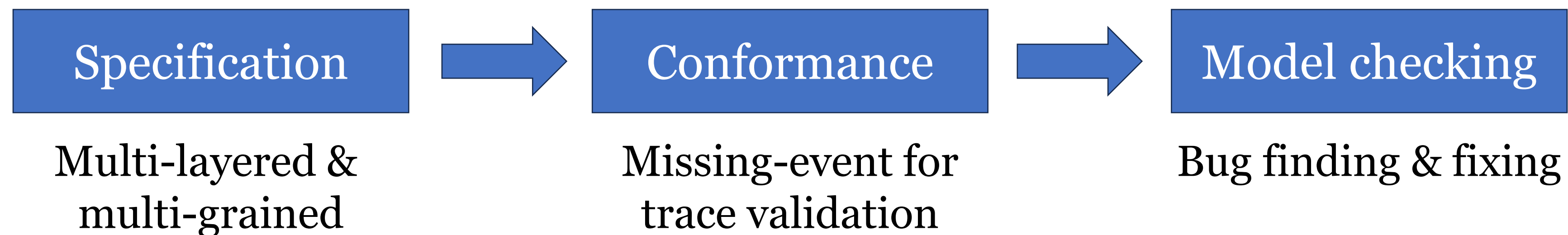
## Challenge 2: Checking spec-code conformance

- **Discrepancies:** Caused by modeling errors or code changes
- **Impact:** Leads to false positives and negatives, reducing confidence

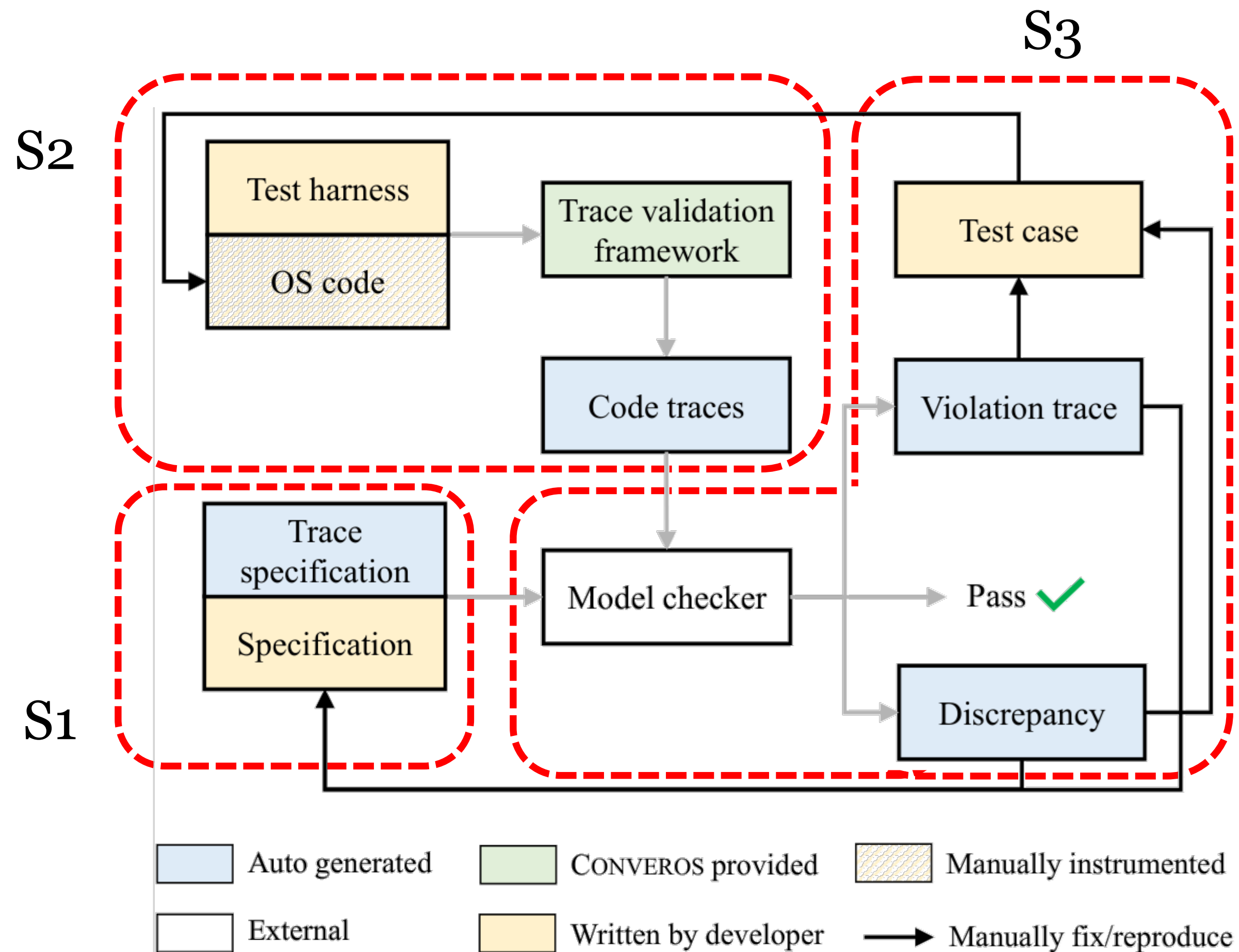


# Contributions

- **CONVEROS: a practical model checking methodology**
  - Workflow: Specification development, conformance checking, model checking
- Applied to **12** concurrency modules in Asterinas OS
  - Found **20** real bugs (e.g., deadlocks, livelocks, data races)
  - Only **4** person-months of effort
  - Low spec-to-code ratio (**0.3–2.3**)
- What makes CONVEROS practical
  - **For specification:** Multi-layered and multi-grained specs to reduce effort and state space
  - **For conformance:** Trace validation with “missing event” support



# CONVEROS overview



- S1: Specification development
  - PlusCal
  - Multi-layered & multi-grained
- S2: Conformance checking
  - Instrumentation for logging
  - Enhanced trace validation (e.g., missing-event support)
  - Revise spec/code for discrepancies
- S3: Model checking
  - Fix code bugs
  - Check the correctness of fixes

# Multi-layered specifications

- Multi-layered: High-level & low-level specs
- **High-level specs**
  - Capture the design intent of sync primitives
  - Model public APIs (e.g., lock(), unlock())
  - Abstract away internal implementation details
  - Benefits
    - Easy to write and verify
    - Enables early conformance checking
    - Serve as a basis for developing low-level specifications

```
1  CONSTANTS PS /* Process set
2  (*--algorithm AbstractLock {
3  variables locked = FALSE;
4  procedure acquire_lock() {
5  lock: /* A label
6    await locked = FALSE;
7    locked := TRUE;
8    return;
9  }
10 procedure release_lock() {
11 unlock:
12   locked := FALSE;
13   return;
14 }
15 fair process (p \in PS) {
16 start:
17   while (TRUE) {
18     call acquire_lock();
19   cs:
20     skip;
21     call release_lock();
22   }
23 }}*)

VARIABLES locked, pc
/* Some variables are omitted
vars == <<locked, pc>>
lock(p) ==
  /\ pc[p] = "lock"
  /\ locked = FALSE
  /\ locked' = TRUE
  /\ pc' = [pc EXCEPT ![p]="cs"]
acquire_lock(p) == lock(p)
start(p) ==
  /\ pc[p] = "start"
  /\ pc' = [pc EXCEPT ![p]="lock"]
  /\ UNCHANGED locked
/* Some actions are omitted
Init ==
  /\ locked = FALSE
  /\ stack = [p \in PS |-> <<>>]
  /\ pc = [p \in PS |-> "start"]
Next == \E p \in PS:
  \V acquire_lock(p)
  \V release_lock(p)
  \V start(p)
  \V cs(p)
```

High-level specification of the lock() and unlock() API  
(written in PlusCal, compiled to TLA+ for verification)

# Multi-layered specifications

## Low-level specs

- Build based on high-level specs to closely reflect implementation
- Use labeled blocks to define atomic actions
- Translate Rust logic into PlusCal processes and actions
  - Capture detailed control logic and state changes
  - Explicitly model Rust-specific semantics (e.g., drop())

```
1 variables locked = FALSE, has_woken_of_wakers = <<>>,
2   num_wakers = 0, wakers_lock = FALSE, wakers = <<>>;
3 procedure acquire_lock()
4   variable waiter_ref; (* Local variable *) {
5   fast_path:
6     if (locked = FALSE) { locked := TRUE; return; };
7   slow_path:
8     has_woken_of_wakers := Append(has_woken_of_wakers, FALSE);
9     waiter_ref := Len(has_woken_of_wakers);
10 loop:
11   while (TRUE) {
12   \*enqueue: skip;
13   acquire_wakers_lock:
14     await wakers_lock = FALSE; wakers_lock := TRUE;
15   enqueue_waker:
16     wakers := Append(wakers, waiter_ref);
17   inc_num_wakers:
18     num_wakers := num_wakers + 1;
19   release_wakers_lock:
20     wakers_lock := FALSE;
21   try_lock:
22     if (locked = FALSE) { locked := TRUE; goto drop_waiter; };
23   do_wait:
24     await has_woken_of_wakers[waiter_ref] = TRUE;
25     has_woken_of_wakers[waiter_ref] := FALSE;
26   };
27 drop_waiter:
28   has_woken_of_wakers[waiter_ref] := TRUE; return; }
```

Low-level spec for Mutex::lock()  
in PlusCal

```
1 let mut guard = mutex.lock();
2 *guard += 1;
3 // `guard` goes out of scope here, lock released via drop()
```

Rust background: implicit unlock() via drop()

```
1 pub struct Waker { has_woken: AtomicBool } // has_woken_of_wakers
2 pub struct Waiter { waker: Arc<Waker> }
3 pub struct WaitQueue { num_wakers: AtomicU32, // num_wakers
4   wakers: SpinLock<VecDeque<Arc<Waker>>>} // wakers_lock, wakers
5 pub struct Mutex<T> {lock: AtomicBool, queue: WaitQueue, val: T}
6 impl<T> Mutex<T> {
7   pub fn lock(&self) -> MutexGuard<T> { // acquire_lock in spec
8     self.queue.wait_until(|| self.try_lock())
9   }
10 impl WaitQueue {
11   pub fn wait_until(&self, cond: ...) -> ... {
12     if let Some(res) = cond() { return res; } // fast_path
13     let (waiter, _) = Waiter::new_pair(); // slow_path
14     let cond = || { self.enqueue(waiter.waker()); // enqueue
15                       cond() }; // try_lock
16     loop { // loop
17       if let Some(res)=cond() { /* drop(waiter); */ return res;};
18       waiter.waker.do_wait(); // do_wait
19     }
20   }
21   pub fn enqueue(&self, waker: Arc<Waker>) {
22     let mut wakers = self.wakers.lock(); // acquire_wakers_lock
23     wakers.push_back(waker); // enqueue_waker
24     self.num_wakers.fetch_add(1); // inc_num_wakers
25     /* drop(wakers); */ // release_wakers_lock
26   }
27 }
28 impl Drop for Waiter { // drop_waiter
29   fn drop(&mut self) { self.waker.has_woken.swap(true); }
```

Corresponding Rust implementation  
of Mutex::lock()

# Multi-grained specifications

- Fine-grained specs lead to state explosion
- Coarse-grained may miss bugs
- We adopt a multi-grained approach
  - Fine-grained specs for basic sync primitives
  - Mix-grained for modules built on top of sync primitives
    - Omit internal details of verified components
    - Coarsen verified behaviors (e.g., model them as atomic variables)
    - Correctness is retained in coarsened specs by the interaction-preserving principle<sup>[1]</sup>

```
3 pub struct WaitQueue { num_wakers: AtomicU32, // num_wakers
4   wakers: SpinLock<VecDeque<Arc<Waker>>>} // wakers_lock, wakers
5 pub struct Mutex<T> {lock: AtomicBool, queue: WaitQueue, val: T}
6 impl<T> Mutex<T> {
7   pub fn lock(&self) -> MutexGuard<T> { // acquire_lock in spec
8     self.queue.wait_until(|| self.try_lock())
9   }}
```

Code snippet of Mutex::lock()

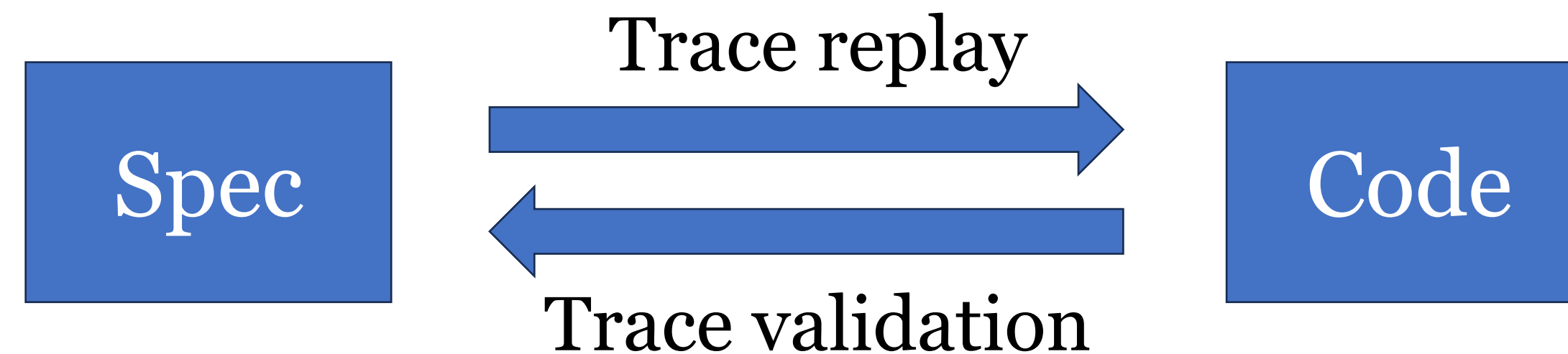
```
1 variables locked = FALSE, has_woken_of_wakers = <<>>,
2   num_wakers = 0, wakers_lock = FALSE, wakers = <<>>;
3 procedure acquire_lock()
4 variable waiter_ref; (* Local variable *) {
5 fast_path:
6   if (locked = FALSE) { locked := TRUE; return; };
7 slow_path:
8   has_woken_of_wakers := Append(has_woken_of_wakers, FALSE);
9   waiter_ref := Len(has_woken_of_wakers);
10 loop:
11   while (TRUE) {
12     \*enqueue: skip;
13     acquire_wakers_lock:
14       await wakers_lock = FALSE; wakers_lock := TRUE;
15     enqueue_waker:
16       wakers := Append(wakers, waiter_ref);
17     inc_num_wakers:
18       num_wakers := num_wakers + 1;
19     release_wakers_lock:
20       wakers_lock := FALSE;
21   try_lock:
22     if (locked = FALSE) { locked := TRUE; goto drop_waiter; };
23   do_wait:
24     await has_woken_of_wakers[waiter_ref] = TRUE;
25     has_woken_of_wakers[waiter_ref] := FALSE;
26   };
27 drop_waiter:
28   has_woken_of_wakers[waiter_ref] := TRUE; return; }
```

Low-level spec for Mutex::lock()

[1] L. Ouyang, X. Sun, R. Tang, Y. Huang, M. Jivrajani, X. Ma, and T. Xu. Multi-Grained Specifications for Distributed System Model Checking and Verification. In Proceedings of the 20th European Conference on Computer Systems (EuroSys), 2025.

# Conformance checking

- **Goal:** Bridge the gap between specification and implementation
- Two common approaches in TLA+
  - Trace replay (spec  $\rightarrow$  code): generate traces from the spec and replay them in the code
    - Challenging for OS-level code: requires deterministic control of thread interleavings
  - **Trace validation** (code  $\rightarrow$  spec): run the code to generate execution traces, and check whether the observed behavior allowed by the specification
    - Practical for OS kernels: only requires lightweight logging instrumentation



# Trace validation

- Trace collection
  - Develop test harnesses and instrument the code according to the spec's labels
  - Trace logs include timestamp, process ID, event name, and variable updates

```
1 CONSTANTS PS \* Process set
2 (--algorithm AbstractLock {
3 variables locked = FALSE;
4 procedure acquire_lock() {
5 lock: \* A label
6   await locked = FALSE;
7   locked := TRUE;
8   return;
9 }
10 procedure release_lock() {
11 unlock:
12   locked := FALSE;
13   return;
14 }
15 fair process (p \in PS) {
16 start:
17   while (TRUE) {
18     call acquire_lock();
19   cs:
20     skip;
21     call release_lock();
22   }
23 }}*)
```

Specification

```
1 pub struct MutexModel;
2 impl TlaModel for MutexModel {
3   fn run(&self, procs: usize, loops: usize) {
4     let mutexlock = Arc::new(Mutex::new(0));
5     let proc_func = move || {
6       for _i in 0..loops {
7         TlaLogger::new("start").next_pc("lock").record();
8         let guard = mutexlock.lock();
9         TlaLogger::new("lock").next_pc("cs").record();
10        TlaLogger::new("cs").next_pc("unlock").record();
11        drop(guard);
12        TlaLogger::new("unlock").record();
13      }
14    };
15    run_procs(procs, proc_func); // Provided by trace framework
16  }
```

Test harness

```
1 {"timestamp":100,"event":"start","p":2,"pc":[{"op":"Update", "path":[2],"args":["lock"]}]}
2 {"timestamp":110,"event":"start","p":1,"pc":[{"op":"Update", "path":[1],"args":["lock"]}]}
```

Trace log

# Trace validation

- Trace specification
  - ① Extends high-level/low-level specs
  - ② Constrains allowed actions
- Validation process
  - Reads trace logs
  - Constrain each log's event and variable updates against the spec
  - Explore allowed (enabled) actions using TLC model checker
  - If a discrepancy is detected
    - At the low-level spec: revise the spec to fix modeling errors
    - At the high-level spec: fix code deviation bugs

```
1  /* Simplified actions of the high-level lock specification
2  start(self) ==
3    /\ pc[self] = "start" /\ pc' = [pc EXCEPT ![self] = "lock"]
4    /\ UNCHANGED locked
5  lock(self) ==
6    /\ pc[self] = "lock" /\ pc' = [pc EXCEPT ![self] = "cs"]
7    /\ locked = FALSE /\ locked' = TRUE
8  /* Trace specification (below)
9  VARIABLES l /* Initialized to 1
10 Trace == LoadTrace("trace.ndjson") /* Loaded from trace file
11 cur_line == Trace[l]
12 UpdateVariables(line) ==
13   "pc" \in DOMAIN line => pc' = UpdateVar(pc, line.pc) /\ ...
14 IsEvent(e) ==
15   /\ l \in 1..Len(Trace) /\ l' = l + 1
16   /\ pc[cur_line.p] = e /* Use PlusCal's `pc` as event constraint
17   /\ UpdateVariables(cur_line)
18 SpecActions(p) == start(p) \/ lock(p) \/ ...
19 TraceNext == IsEvent(cur_line.event) /\ SpecActions(cur_line.p)
```

Trace spec

```
1 {"timestamp":100,"event":"start","p":2,"pc":[{"op":"Update", "path":[2],"args":["lock"]}]}
2 {"timestamp":110,"event":"start","p":1,"pc":[{"op":"Update", "path":[1],"args":["lock"]}]}

```

Trace log

# Enhanced trace validation

- Issues:
  - Logged timestamps may be inaccurate
  - Logging every label can be burdensome
- Our improvement: Allow missing events
  - Let TLC infer valid events
  - Insert “missing event” logs before and after an action to define a possible time window
  - Consecutive missing events are merged
    - Ensures that each missing event is followed by a normal event

```
1 pub struct MutexModel;
2 impl TlaModel for MutexModel {
3   fn run(&self, procs: usize, loops: usize) {
4     let mutexlock = Arc::new(Mutex::new(0));
5     let proc_func = move || {
6       for _i in 0..loops {
7         TlaLogger::new("start").next_pc("lock").record();
8         TlaLogger::new_missing().cur_pc("lock").record();
9         let guard = mutexlock.lock();
10        TlaLogger::new_missing().cur_pc("lock").record();
11        TlaLogger::new("cs").next_pc("unlock").record();
12        TlaLogger::new_missing().cur_pc("unlock").record();
13        drop(guard);
14        TlaLogger::new_missing().cur_pc("unlock").record();
15      }
16    };
17    run_procs(procs, proc_func);
18  }
19 }
```

Test harness (missing event)

```
1 {"timestamp":110,"event":"MissingEvent","ps":{"1,2},"cur_pc":[{"op":"in","path":[1],"args":["lock"]}, {"op":"in","path":[2],"args":["lock"]}]}
2 {"timestamp":120,"event":"cs","p":1}
3 {"timestamp":130,"event":"MissingEvent","ps":{"1,2},"cur_pc":[{"op":"in","path":[1],"args":["lock"]}, {"op":"in","path":[2],"args":["lock"]}]}
4 {"timestamp":140,"event":"cs","p":2,"converged":[1]}
```

Trace log (missing event)

# Enhanced trace specification

- Handling missing events in trace specification
  - **Branch ①**: Matches the next non-missing event
  - **Branch ②**: Explores all enabled actions without consuming the missing event, keeping it active for further steps
- Limit search space from missing events for efficiency
  - ① Max depth pruning (e.g., depth < 10)
  - ② Convergence points: stop exploration for a missing event once its corresponding convergence log is encountered

```

1 next_line == Trace[l + 1]
2 MissingEvent(Actions(_)) == \* Actions: operator with 1 argument
3 ① /\ depth < cur_line.max_depth /\ depth' = depth + 1 \* Max depth pruning
4 ② /\ TLGGet(1) \* Converge point pruning
5   /\ UNCHANGED 1 \* ps (below): allowed process set
6   /\ \E p \in cur_line.ps: UpdateVariables(cur_line) /\ Actions(p)
7 IsNextEvent(e) ==
8   /\ l \in 1..(Len(Trace)-1) /\ l' = l + 2
9   /\ pc[next_line.p] = e
10  /\ UpdateVariables(next_line)
11 ② /\ \A i \in cur_line.converged: TLCSet(i, FALSE) \* Set converge points
12 ExecuteEvent(Actions(_)) ==
13   IF cur_line.event /= "MissingEvent"
14   THEN IsEvent(cur_line.event) /\ Actions(cur_line.p)
15 ① ELSE \/\ IsNextEvent(next_line.event) /\ Actions(next_line.p)
16 ② \/\ MissingEvent(Actions)
17 TraceNext == ExecuteEvent(SpecActions)

```

Trace spec (missing event)

```

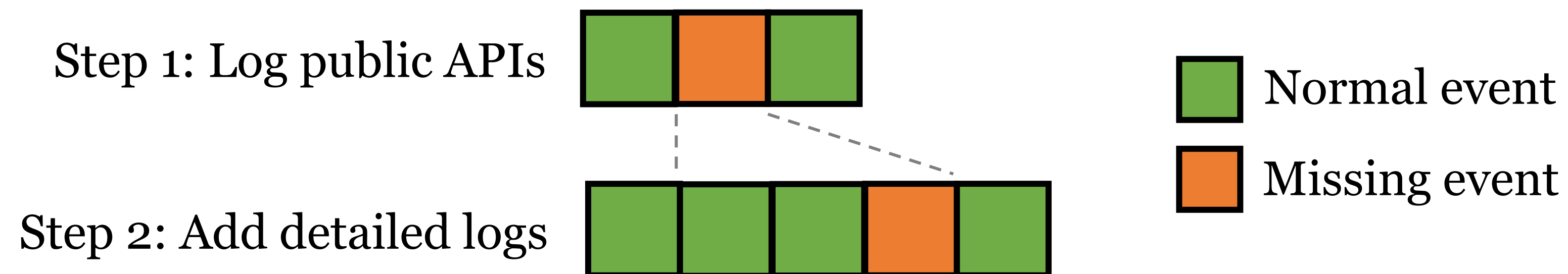
1 {"timestamp":110,"event":"MissingEvent","ps":{"1,2},"cur_pc":[{"op":"in","path":[1],"args":["lock"]}, {"op":"in","path":[2],"args":["lock"]}]}
2 {"timestamp":120,"event":"cs","p":1}
3 {"timestamp":130,"event":"MissingEvent","ps":{"1,2},"cur_pc":[{"op":"in","path":[1],"args":["lock"]}, {"op":"in","path":[2],"args":["lock"]}]}
4 {"timestamp":140,"event":"cs","p":2,"converged":[1]}

```

Trace log (missing event)

# Enhanced trace collection

- Incremental trace collection
  - Start from public APIs of synchronization primitives
  - Gradually instrument internal statements
  - Enabled by missing events
- Benefits
  - Enables early detection and debugging of spec-code discrepancies
  - Reduces logging overhead in early development stages



# Model checking

- Use TLC model checker to check for property violations
  - E.g., common safety and liveness properties

```
1 MutualExclusion      == \* Safety
2   \A i, j \in Procs: (i /= j) => ~(pc[i] = "cs" /\ pc[j]="cs")
3 DeadAndLiveLockFree == \* Liveness
4   \E i \in Procs: pc[i]="lock" ~> (\E j \in Procs: pc[j]="cs")
5 StarvationFree      == \* Liveness
6   \A i \in Procs: pc[i]="lock" ~> (pc[i]="cs")
```

- Create test cases based on violation traces
  - If the bug is reproducible
    - Fix the spec and the code
    - Rerun trace validation to check conformance
  - If the bug is not reproducible
    - The discrepancy likely lies in the spec, revise the spec
- By-product bugs (e.g., panics or hangs) may also be uncovered during test execution

# Evaluation

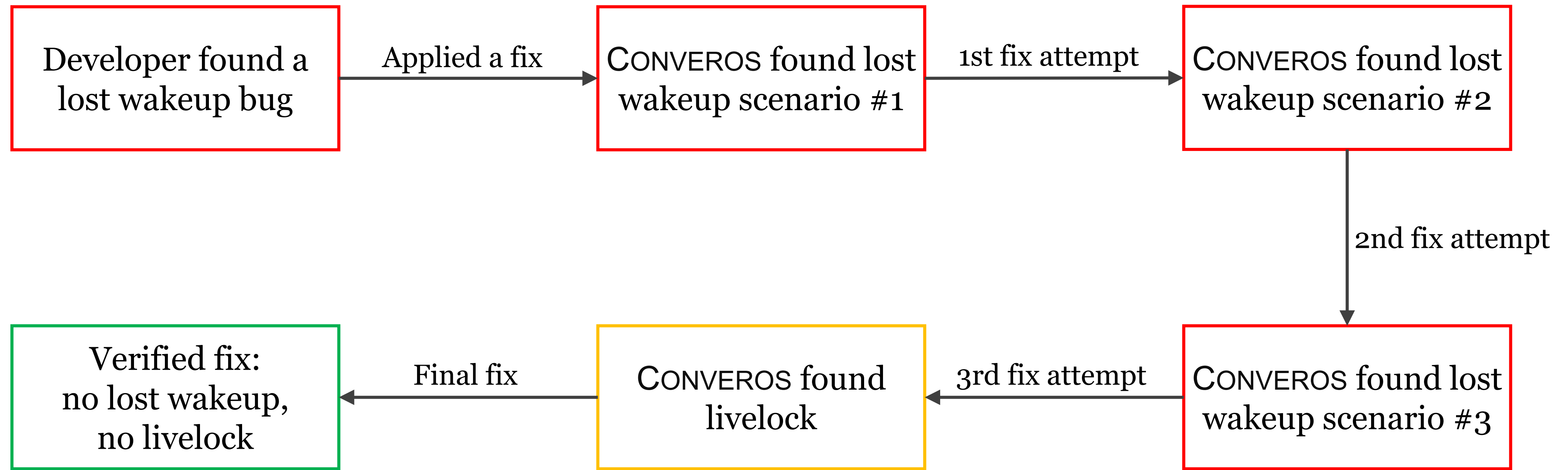
- Verified **12** concurrency modules in Asterinas OS
  - ~4,000 lines of Rust code
- Bugs found: **20** in total
  - 9 from model checking
  - 11 by-product bugs
    - 2 found during modeling
    - 9 found during test case execution
  - All confirmed, 14 fixed
- Effort: **~4 person-months** in total

# Bugs found

ID	Module	Violation/Manifestation	Description
<i>Found by CONVEROS</i>			
#1	RangeLock	Deadlock-free	Lost wakeup due to waiting on an outdated wait queue
#2	Mutex	Mutual exclusion	Incorrect guard construction in lock() causing unintended unlock() on drop()
#3	RwLock	Livelock/Starvation-free	Non-atomic read() operation racing with downgrade()
#4	RwMutex	Deadlock-free	Wrong wakeup condition causing a dropped upreader to never wake up a waiter
#5	Semaphore	Semaphore count $\geq 0$	Semaphore count decreases below 0 due to unhandled multi-operation semop()
#6	Futex	Deadlock-free	Signal or timeout leaving outdated futex item uncleared, causing lost wakeup
#7	Flock	Deadlock-free	Incorrect wait condition causing permanent wait
#8	Pipe	Write atomicity	Writes smaller than PIPE_BUF size not guaranteed to be atomic
#9	TTY	Deadlock-free	Circular dependencies from incorrect locking order causing TTY hang
<i>By-product bugs</i>			
#10	Semaphore	(Found during modeling)	Incorrect initialization of the waiter count, leading to wrong lookup results
#11	Futex	(Found during modeling)	Poor hash algorithm leads to uneven key distribution, overloading a few buckets
#12	Atomic mode	Kernel panic	Misuse of spinlock/mutex breaks atomic mode
#13	IRQ	Stack overflow	Improper IRQ enabling during the bottom half leads to nested IRQs
#14	TLB	Kernel hang	Deadlock caused by a spinlock used in IRQ without disabling IRQ outside it
#15	Trap	Kernel panic and hang	Uncleared direction flag (rflags.DF) leads to memory corruption on traps
#16	Task switch	User process crash	Missing FPU state save/restore causes memory corruption on task switch
#17	Sendfile syscall	User data truncation	Short writes cause data truncation when input is read but not fully written
#18	Kernel tests	Kernel test hang	Assertion in kernel test threads leads to infinite error message unwraps
#19	Kernel threads	Kernel panic	New scheduler feature breaks CPU affinity setting
#20	Logging	Kernel hang	Logging allocating memory during low-memory rescue triggers recursive rescues

# Example bug: Hard-to-fix lost wakeup in RangeLock

- RangeLock's (i.e., file record lock) complex semantics make it easy to introduce subtle concurrency bugs
- Developer found and fixed a lost wakeup bug, but the fix was flawed
- CONVEROS helped identify issues in each fix until the bug was correctly resolved



# Example discrepancy: Subtle modeling error in RwMutex

- Discrepancy in try\_upread logic
  - Code: sets the UPREADER bit and checks previous lock state (WRITER | UPREADER)
  - Buggy spec: checked current writer\_lock after update (not previous state)
  - Correct spec: stores previous states of both writer\_lock and upreader\_lock
  - Found by trace validation: writer\_lock changed before the check label

```
1 pub fn try_upread(&self) -> Option<RwMutexUpReadGuard<T>> {  
2   let lock = self.lock.fetch_or(UPREADER) & (WRITER | UPREADER);  
3   if lock == 0 {  
4     Some(RwMutexUpReadGuard { inner: self })  
5   }
```

RwMutex implementation

```
1 procedure try_upread()  
2 variable prev_lock; {  
3   fetch_or:  
4   - prev_lock := upreader_lock;  
5   + prev_lock := <<upreader_lock, writer_lock>>;  
6   upreader_lock := TRUE;  
7   check: ←  
8   - if (~prev_lock /\ ~writer_lock) {  
9   + if (~prev_lock[1] /\ ~prev_lock[2]) {  
10    success:  
11    role[self] := UPREADER;  
12    return;  
13  }
```

RwMutex specification (discrepancy and fix)

Another thread holding the write lock released it between the “fetch\_or” and “check” labels

# Verification effort

Code		Specification			Estimate effort	
Module	#LOC	#LOC	#Var.	#Act.	Spec.	Conf.
SpinLock	147	49	1	5	1.5	5
Mutex	113	89	4	14	4	3
RwLock	404	261	5	19	5	2.5
RwMutex	202	460	9	47	3.5	2.5
CondVar	238	199	6	25	3.5	1.5
Semaphore	622	490	10	37	9	2
PageCursor	468	297	10	21	5	-
Pipe	392	143	7	17	1.5	1
RangeLock	457	529	9	24	4	2
Flock	144	176	7	12	1.5	0.5
Futex	469	201	5	21	4	1
TTY	309	138	4	20	1	0.5
<b>Total</b>	<b>3965</b>	<b>3032</b>	<b>77</b>	<b>262</b>	<b>43.5</b>	<b>21.5</b>

- Specification-to-code ratio
  - Ranges from 0.3 to 2.3
- Specifications development
  - 43.5 person-days
- Conformance checking
  - 21.5 person-days
- Total effort
  - ~4 person-months

# Conclusion

- CONVEROS: a practical model checking methodology
  - 3-step workflow: specification → conformance checking → model checking
  - Multi-layered, multi-grained specifications for usability and scalability
  - Enhanced trace validation with missing-event support for trace validation
- Evaluation summary
  - Verified **12** concurrency modules (~4,000 LoC) in Asterinas OS
  - Found **20** bugs (e.g., deadlocks, livelocks, kernel panics)
  - ~4 person-months of effort; spec-to-code ratio: 0.3–2.3

**THANK YOU!**