Inherently Deterministic Operating Systems

Fachgruppe Betriebssysteme, Herbsttreffen 2021, 2021-09-22

Stefan Reif¹, Timo Hönig², Wolfgang Schröder-Preikschat¹

¹Friedrich-Alexander-Universität Erlangen-Nürnberg ²Ruhr-Universität Bochum







Scalable systems

- $\blacksquare>$ 100 cores
- communicating partitions
- () system noise
- () tail latencies

Scalable systems

- > 100 cores
- communicating partitions
- I system noise
- () tail latencies

Predictable systems

- pprox 1 core (isolated partitions)
- () low performance
- bad energy efficiency
- 1 poor utilisation



Interferences - Toy Example



- OSs orchestrate communication
 - Interferences due to shared data structures
- OSs manage shared stateful data
 - F Need for synchronisation
- How to build operating systems that are ...
 - ⑦ scalable
 - ⑦ predictable
 - ⑦ and analysable?

Approach

- Inherently Deterministic Operating Systems ("IDOS")
 - \rightarrow IDOS := costs are independent of...
 - ...system states & input data
 - ...the number of resources
 - ...concurrent operations
- Predictability & concurrency as afterthoughts first-class citizens

- Inherently Deterministic Operating Systems ("IDOS")
 - \rightarrow IDOS := costs are independent of...
 - ...system states & input data \Rightarrow single-path code
 - ...the number of resources $\Rightarrow \mathcal{O}(1)$ algorithms
 - ...concurrent operations \Rightarrow wait-freedom
- Predictability & concurrency as afterthoughts first-class citizens

- Inherently Deterministic Operating Systems ("IDOS")
 - \rightarrow IDOS := costs are independent of...
 - ...system states & input data \Rightarrow single-path code
 - ...the number of resources $\Rightarrow \mathcal{O}(\mathbf{1})$ algorithms if possible
 - ...concurrent operations \Rightarrow wait-freedom but no universal $\mathcal{O}(1)$ approach
- Predictability & concurrency as afterthoughts first-class citizens

- Inherently Deterministic Operating Systems ("IDOS")
 - \rightarrow IDOS := costs are independent of...
 - ...system states & input data \Rightarrow single-path code
 - ...the number of resources $\Rightarrow \mathcal{O}\left(\textbf{1} \right)$ algorithms if possible
 - ...concurrent operations \Rightarrow wait-freedom but no universal $\mathcal{O}(1)$ approach
- Predictability & concurrency as afterthoughts first-class citizens
 - \odot Scalability: $\mathcal{O}(1) < \mathcal{O}(P)$ for large P
 - Predictability: costs are constant
 - Analysability: trivial, no overestimation, incremental
 - ⑦ Implementability?
 - ⑦ Efficiency?

Implementation of an IDOS prototype

- DETOX—the Deterministic Operating System
- Identify ...
 - 1. ...hardware requirements
 - () Constant-time multi-core hardware does not exist
 - \Rightarrow {Functional, non-functional} hardware features?
 - \Rightarrow {Necessary, desirable} hardware features?
 - 2. ...suitable algorithms & data structures
 - Non-universality
 - \Rightarrow Concurrent $\mathcal{O}(1)$ operations
 - 3. ...implementable application-level interface
 - \Rightarrow Deterministic coordination & communication

- 1. Approach
- 2. Design
- 3. Implementation
- 4. Constant-Time Programming
- 5. Discussion
- 6. Evaluation
- 7. Conclusion

Design

Design of DETOX | Strategies

System behaviour

- 1. Explicit thread placement
 - Deterministic control flow location
 - Thread migration supported
- 2. Cooperative scheduling
 - No involuntary preemptions
 - Voluntary release of CPU supported
- 3. Analysable memory demand
 - Constant memory demand per object
 - No need for run-time allocation
- 4. No resource limitations
 - Dynamic object creation supported

Design of DETOX | Strategies

System behaviour

- 1. Explicit thread placement
 - Deterministic control flow location
 - Thread migration supported
- 2. Cooperative scheduling
 - No involuntary preemptions
 - Voluntary release of CPU supported
- 3. Analysable memory demand
 - Constant memory demand per object
 - No need for run-time allocation
- 4. No resource limitations
 - Dynamic object creation supported

Design guidelines

- Worst-case minimisation
 - Worst-case = average-case
- Data access restrictions
 - Data-structure partitioning
 - Implicit ownership model
- Redundancy elimination
 - Avoid inconsistency correction overhead
 - Consider concurrent updates
- Corner case elimination
 - Invariants for Fail-safe operations
 - Unconditional branch execution

Synchronising queues

- Communication & coordination
- Efficient wait-free signalling

Core management threads

- Thread state changes as events
- Synchronisation simplification
- Migration-based synchronisation
 - Mutual exclusion by thread migration

Companion cubes

- Fail-safe concurrent memory management
- Support for preallocation
- DETOX semaphor interface
 - Thread interaction unification
 - Corner case elimination
- Concurrent ownership model
 - Restrict & control data access

Implementation

- Dynamic memory management
 - memory alloc, free
- Dynamic thread management
 - thread create, exit, join, ...
- Actor-based synchronisation
 - actor create, run
 - (SPSC) reply create, set, get, ...

- Message-based coordination
 - (MPSC) pipe create, recv, send, ...
- Migration-based synchronisation
 - thread migrate, pause
- Lock-based synchronisation
 - mutex create, enter, leave
 - condition create, wait, wake

Constant-Time Programming

High-level Algorithm

- 1. Start with wait-free code (with $\mathcal{O}(1)$ complexity)
- 2. Transform loops
 - Discussed in literature
 - Practical solution: unrolling
 - Little relevance due to $\mathcal{O}(\mathbf{1})$ code
- 3. Transform conditions
 - Discussed in literature
 - Problem: atomic memory operations
 - "interim & choose"
- 4. Avoid run-time errors
 - Division by zero, invalid pointer dereference, ...
 - "interim & choose" again







Stefan Reif et al.



Selection Implementation

Stefan Reif et al.





Code Transformation for IDOS | Example

- Address calculation adapted to consider condition
- Unconditional execution of functional part
- Code more difficult to write & understand

Code Transformation for IDOS | Example

```
void enq_if(bool cond,
        queue_t *list, chain_t *item)
{
     f(cond) {
        chain_t **addr = &list->tail; c
        chain_t *prev = XCHG(addr, item);
        prev->next = item;
     }
     p
}
```

- Address calculation adapted to consider condition
- Unconditional execution of functional part
- Code more difficult to write & understand

Code Transformation for IDOS | Example

```
void eng if(bool cond,
                                            void eng if(bool cond,
           gueue t *list, chain t *item)
                                                        gueue t *list. chain t *item)
 if (cond) {
                                             chain t temp = { .next = &temp };
   chain t **addr = &list->tail;
                                             chain t **addr = cond ? &list->tail
   chain t *prev = XCHG(addr. item);
                                             chain t *prev = XCHG(addr, item);
   prev->next = item:
                                             prev->next = item;
```

- Address calculation adapted to consider condition
- Unconditional execution of functional part
- Code more difficult to write & understand

: &temp.next:

Discussion

Discussion of DETOX

- Efficiency
 - Worst-case costs apply unconditionally
 - No overly pessimistic analysis needed
- Implementation issues
 - Device drivers, energy efficiency (Hardware-software interface)
 - Unimplemented features (Memory Protection, ...)

Discussion of DETOX

- Efficiency
 - Worst-case costs apply unconditionally
 - No overly pessimistic analysis needed
- Implementation issues
 - Device drivers, energy efficiency (Hardware-software interface)
 - Unimplemented features (Memory Protection, ...)
- Fundamental issues
 - Functional timing dependencies
 - Communication (e.g., reply.get)
 - Blocking operations (e.g., thread.join)
 - Approach is not universal

Discussion of DETOX

- Efficiency
 - Worst-case costs apply unconditionally
 - No overly pessimistic analysis needed
- Implementation issues
 - Device drivers, energy efficiency (Hardware-software interface)
 - Unimplemented features (Memory Protection, ...)
- Fundamental issues
 - Functional timing dependencies
 - Communication (e.g., reply.get)
 - Blocking operations (e.g., thread.join)
 - Approach is not universal
- Time to say goodbye to...
 - ...interrupts, traps & preemption
 - ...non- $\mathcal{O}(1)$ scheduling strategies
 - ...

Co-existence of Nondeterministic Subsystems



Evaluation

Control-Flow Graphs

reply_set:		
mov	0x8(%rdi),%rdx	
mov	\$0x1,%eax	
mov	%rsi,(%rdi)	
movq	\$0x0,-0x18(%rsp)	
xchg	%rax,(%rdx)	
xor	%esi,%esi	
lea	-0x28(%rsp),%rcx	
test	%rax,%rax	
lea	-0x48(%rsp),%rdi	
sete	%sil	
movl	\$0x0,-0x3c(%rsp)	

sup	Yray Yrcy
Sub	/or dr, /or Cr
sub	%rdx,%rdi
imul	%rsi,%rcx
imul	%rsi,%rdi
mov	0x10(%rcx,%rax,1),%r8
lea	-0x3c(%rsp),%rax
lea	-0x38(%rsp),%rcx
sub	%rdx,%rcx
sub	%r8,%rax
imul	%rsi,%rax
imul	%rsi,%rcx
movslq	(%rax,%r8,1),%rax

mov	%r8,0x10(%rdi,%rdx,1)
movq	\$0x0,(%rcx,%rdx,1)
lea	-0x30(%rsp),%rcx
shl	\$0x6,%rax
mov	%rcx,-0x30(%rsp)
add	\$0x603120,%rax
sub	%rax,%rcx
imul	%rsi,%rcx
mov	%rdx,%rsi
xchg	%rsi,(%rcx,%rax,1)
mov	%rdx,(%rsi)
retq	

\Rightarrow No branches, no control-flow dependency on system state or size

Stefan Reif et al.

Control-Flow Analysis



 \Rightarrow Software-level analysis is almost trivial, only in-hardware timing variance remains

Stefan Reif et al.

Hardware Predictability



- Cycle-accurate RISC-V emulator
- Rocket-chip configuration
 - 2 "small" cores
 - L1D\$.nsets = 2
 - L1D\$.nways = 1
 - L2D\$.nsets = 2
 - L2D\$.nways = 4
 - L2D\$.replacement: TrueLRU
 - Bus arbitration: LowestIndexFirst
 - Cycle measurement at core # o
 - Measurement overhead = 1 cy

Calls with debug interrupt ignored

\Rightarrow Some observable patterns in latency variation, some noise

Scalability



 \Rightarrow Latency independent of parallelism, except hyper-threading (amd64)

Stefan Reif et al.

Conclusion

Summary & Conclusion

- Need for scalable, predictable, and analysable systems
 - Main issue: transitive interferences
 - System structure: interacting subsystems
 - Focus on predictable communication
- Approach: inherently deterministic operating systems (IDOSs)
 - Scalability: $\mathcal{O}(1)$
 - Analysability: almost trivial, no overestimation, incremental
 - ✓ Predictability: constant costs
 - 🔉 ...above ISA-level

Summary & Conclusion

- Need for scalable, predictable, and analysable systems
 - Main issue: transitive interferences
 - System structure: interacting subsystems
 - Focus on predictable communication
- Approach: inherently deterministic operating systems (IDOSs)
 - Scalability: $\mathcal{O}(1)$
 - Analysability: almost trivial, no overestimation, incremental
 - ✓ Predictability: constant costs
 - 🔉 ...above ISA-level

- Co-existence with non-deterministic subsystems
 - Various means for coordination
 - Constraints for functional timing dependencies
- Determinism still has some open issues
 - 🗴 constant-time hardware
 - 🔉 Fundamental limitations (e.g., interrupts)
 - ⑦ Hardware-software interfaces
 - ? Suitable $\mathcal{O}(1)$ strategies (e.g., scheduler)
 - ? Concurrent $\mathcal{O}(1)$ algorithms