



# OPERATING SYSTEMS MEET FAULT TOLERANCE

## Microkernel-Based Operating Systems

Björn Döbel

Dresden, 29.01.2013

*“If there’s more than one possible outcome of a job or task, and one of those outcome will result in disaster or an undesirable consequence, then somebody will do it that way.” (Edward Murphy jr.)*

# Outline

- Murphy and the OS: Is it really that bad?
- Fault-Tolerant Operating Systems
  - Minix3
  - CuriOS
  - L4ReAnimator
- Creative OS Debugging
  - Detecting race conditions with DataCollider

# Why Things go Wrong

- Programming in C:

*This pointer is certainly never going to be NULL!*

# Why Things go Wrong

- Programming in C:

*This pointer is certainly never going to be NULL!*

- Layering vs. responsibility:

*Of course, someone in the higher layers will already have checked this return value.*

# Why Things go Wrong

- Programming in C:

*This pointer is certainly never going to be NULL!*

- Layering vs. responsibility:

*Of course, someone in the higher layers will already have checked this return value.*

- Concurrency:

*This struct is shared between an IRQ handler and a kernel thread. But they will never execute in parallel.*

# Why Things go Wrong

- Programming in C:

*This pointer is certainly never going to be NULL!*

- Layering vs. responsibility:

*Of course, someone in the higher layers will already have checked this return value.*

- Concurrency:

*This struct is shared between an IRQ handler and a kernel thread. But they will never execute in parallel.*

- Hardware interaction:

*But the device spec said, this was not allowed to happen!*

# Why Things go Wrong

- Programming in C:

*This pointer is certainly never going to be NULL!*

- Layering vs. responsibility:

*Of course, someone in the higher layers will already have checked this return value.*

- Concurrency:

*This struct is shared between an IRQ handler and a kernel thread. But they will never execute in parallel.*

- Hardware interaction:

*But the device spec said, this was not allowed to happen!*

- Hypocrisy:

*I'm a cool OS hacker. I won't make mistakes, so I don't need to test my code!*

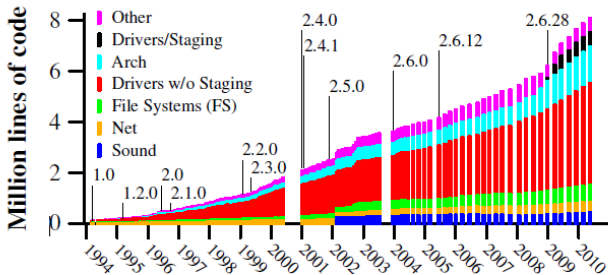
## A Classic Study

- A. Chou et al.: *An empirical study of operating system errors*, SOSP 2001
- Automated software error detection (today:  
<http://www.coverity.com>)
- Target: Linux (1.0 - 2.4)
  - Where are the errors?
  - How are they distributed?
  - How long do they survive?
  - Do bugs cluster in certain locations?

## Revalidation of Chou's Results

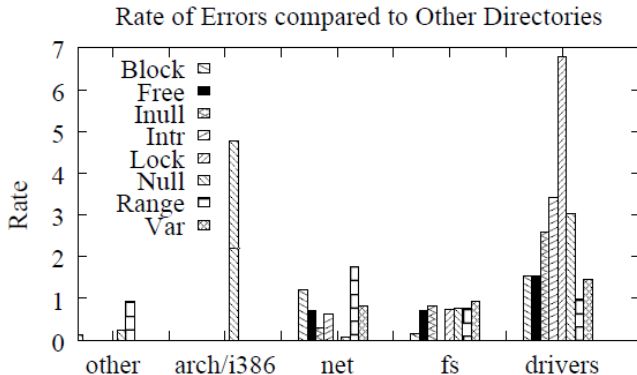
- N. Palix et al.: *Faults in Linux: Ten years later*, ASPLOS 2011
- 10 years of work on tools to decrease error counts - has it worked?
- Repeated Chou's analysis until Linux 2.6.34

# Linux: Lines of Code

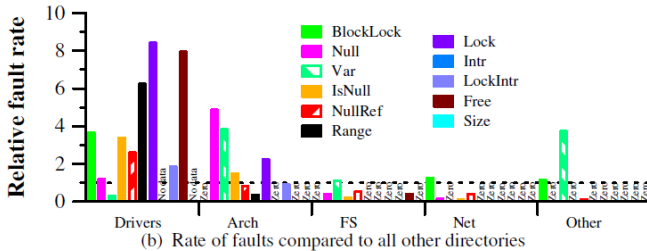


**Figure 1.** Linux directory sizes (in MLOC)

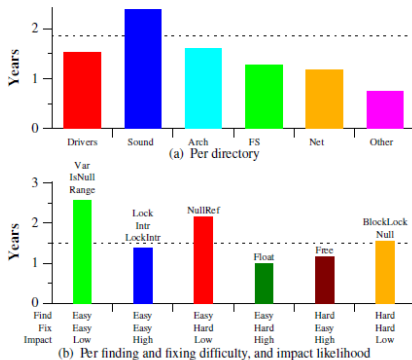
## Fault Rate per Subdirectory (2001)



## Fault Rate per Subdirectory (2011)



## Bug Lifetimes (2011)



**Figure 13.** Average fault lifespans (without staging)

# Bug Distribution

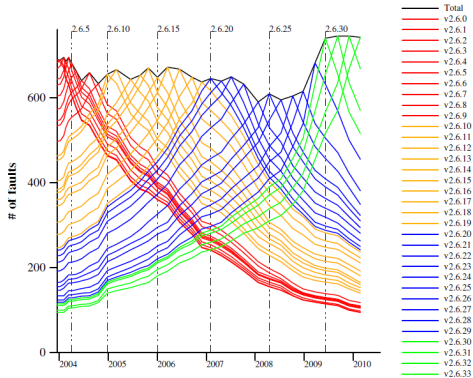
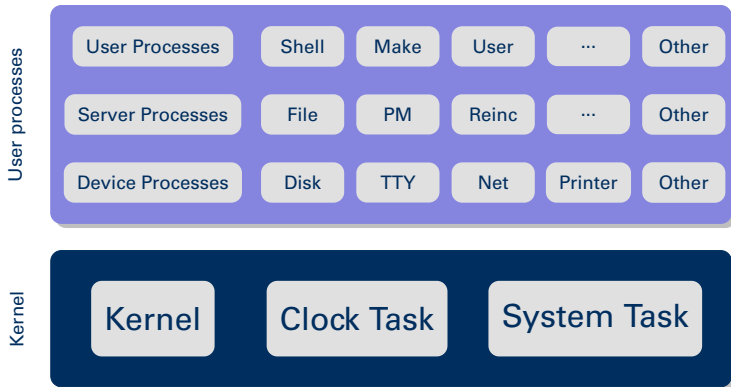


Figure 16. Lifetime of faults across versions

# Break

- Faults are an issue.
- Hardware-related stuff is worst.
- Now what can the OS do about it?

## Minix3 – A Fault-tolerant OS



## Minix3: Fault Tolerance

- Address Space Isolation
  - Applications only access private memory
  - Faults do not spread to other components
- User-level OS services
  - Principle of Least Privilege
  - Fine-grain control over resource access
    - e.g., DMA only for specific drivers
- Small components
  - Easy to replace (micro-reboot)

## Minix3: Fault Detection

- Fault model: transient errors caused by software bugs
- Fix: Component restart
- *Reincarnation server* monitors components
  - Program termination (crash)
  - CPU exception (div by 0)
  - Heartbeat messages
- Users may also indicate that something is wrong

# Repair

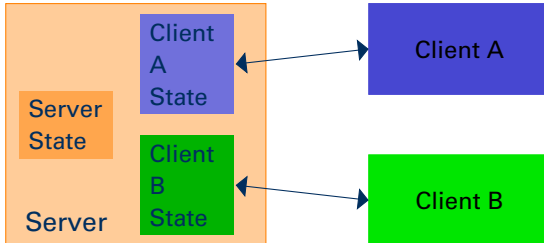
- Restarting a component is insufficient:
  - Applications may *depend* on restarted component
  - After restart, *component state* is lost
- Minix3: explicit mechanisms
  - Reincarnation server signals applications about restart
  - Applications store state at data store server
  - In any case: program interaction needed
    - Restarted app: store/recover state
    - User apps: recover server connection

# Break

- Minix3 fault tolerance:
  - Architectural Isolation
  - Explicit monitoring and notifications
- Other approaches:
  - CuriOS: smart session state handling
  - L4ReAnimator: semi-transparent restart in a capability-based system

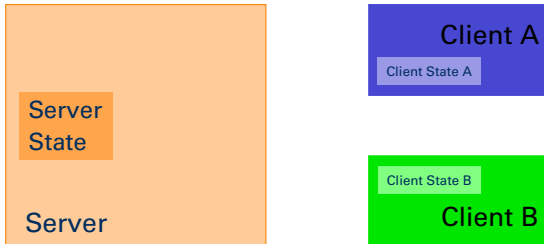
## CuriOS: Servers and Sessions

- State recovery is tricky
  - Minix3: Data Store for application data
  - But: applications interact
    - Servers store *session-specific* state
    - Server restart requires potential rollback for every participant



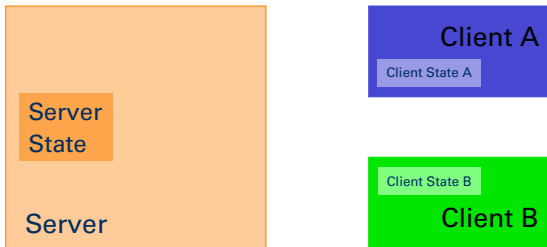
## CuriOS: Server State Regions

- CuriOS kernel (CuiK) manages dedicated session memory: *Server State Regions*
- SSRs are managed by the kernel and attached to a client-server connection



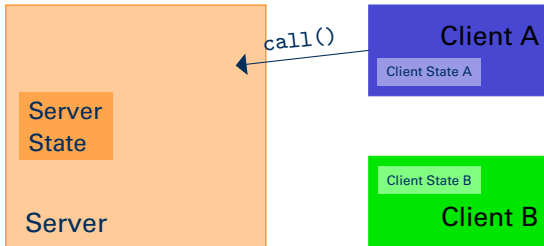
## CuriOS: Protecting Sessions

- SSR gets mapped only when a client actually invokes the server
- Solves another problem: failure while handling A's request will never corrupt B's session state



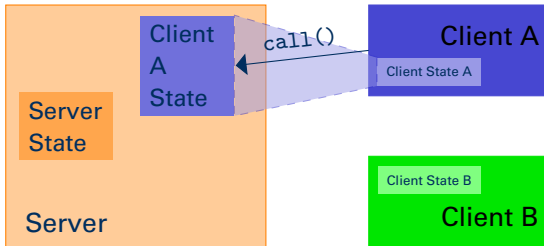
## CuriOS: Protecting Sessions

- SSR gets mapped only when a client actually invokes the server
- Solves another problem: failure while handling A's request will never corrupt B's session state



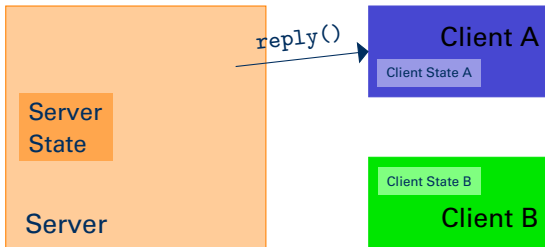
## CuriOS: Protecting Sessions

- SSR gets mapped only when a client actually invokes the server
- Solves another problem: failure while handling A's request will never corrupt B's session state



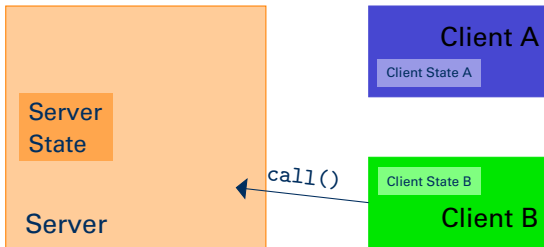
## CuriOS: Protecting Sessions

- SSR gets mapped only when a client actually invokes the server
- Solves another problem: failure while handling A's request will never corrupt B's session state



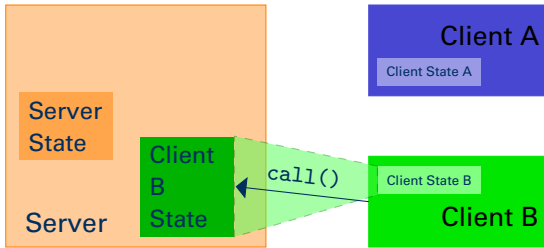
## CuriOS: Protecting Sessions

- SSR gets mapped only when a client actually invokes the server
- Solves another problem: failure while handling A's request will never corrupt B's session state



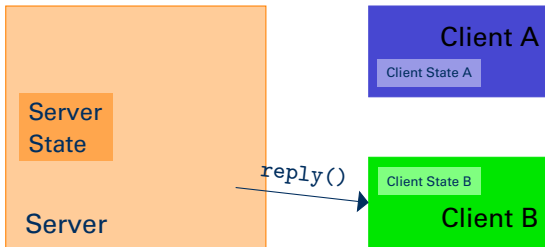
## CuriOS: Protecting Sessions

- SSR gets mapped only when a client actually invokes the server
- Solves another problem: failure while handling A's request will never corrupt B's session state



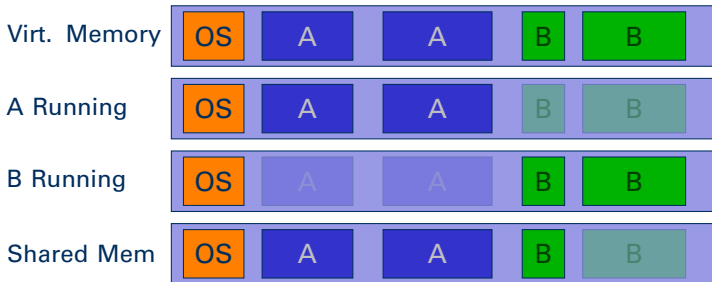
## CuriOS: Protecting Sessions

- SSR gets mapped only when a client actually invokes the server
- Solves another problem: failure while handling A's request will never corrupt B's session state



## CuriOS: Transparent Restart

- CuriOS is a *Single-Address-Space OS*:
  - Every application runs on the same page table (with modified access rights)



# Transparent Restart

- Single Address Space
  - Each object has unique address
  - Identical in all programs
  - Server := C++ object
- Restart
  - Replace old C++ object with new one
  - Reuse previous memory location
  - References in other applications remain valid
  - OS blocks access during restart

## L4ReAnimator: Restart on L4Re

- L4Re Applications
  - Loader component: ned
  - Detects application termination: parent signal
  - Restart: re-execute Lua init script (or parts of it)
  - Problem after restart: capabilities
    - No single component knows everyone owning a capability to an object
    - Minix3 signals won't work

## L4Re: Session Creation

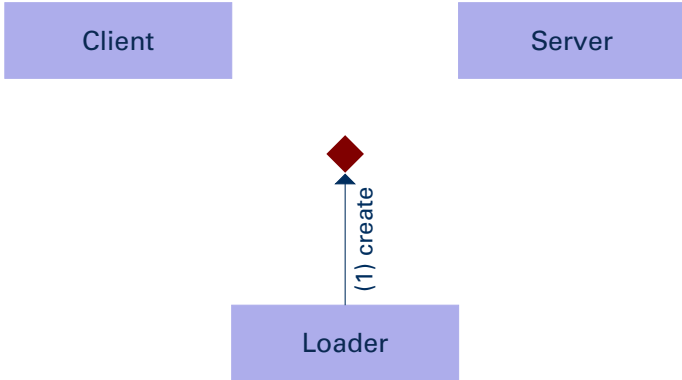
Client

Server

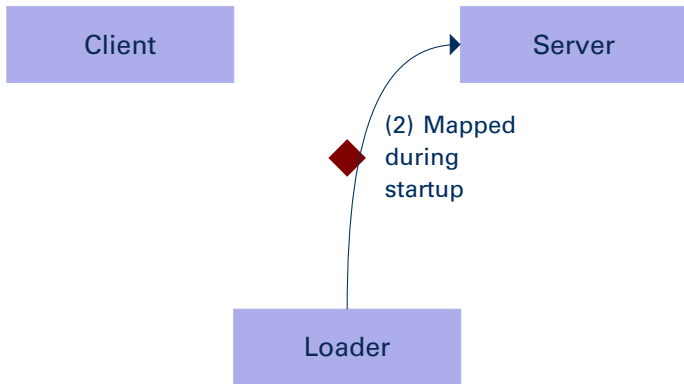
Session  
Creation  
Capability 

Loader

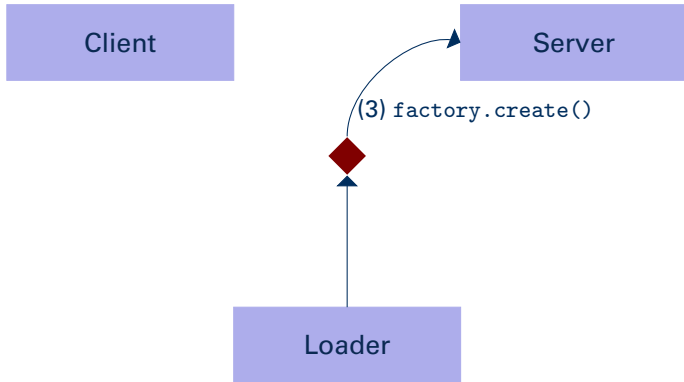
## L4Re: Session Creation



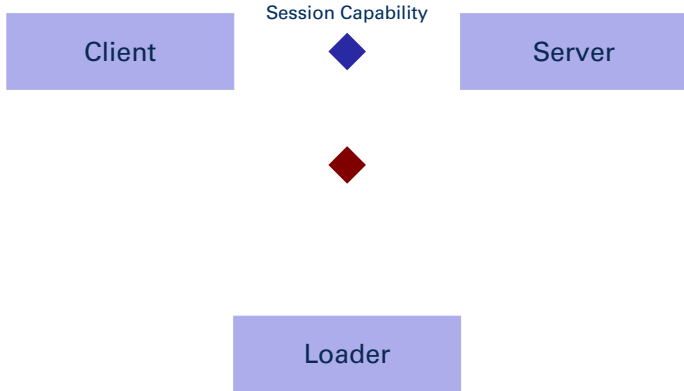
## L4Re: Session Creation



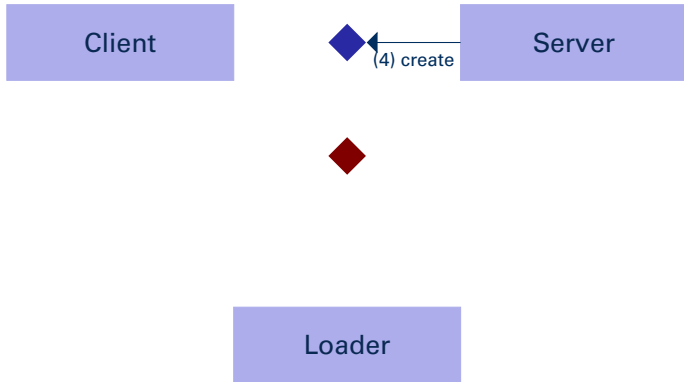
## L4Re: Session Creation



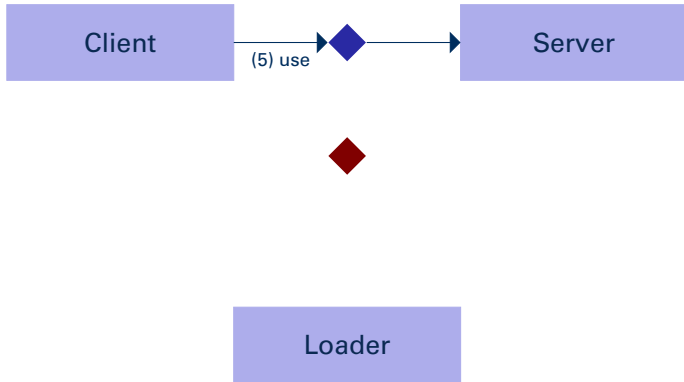
## L4Re: Session Creation



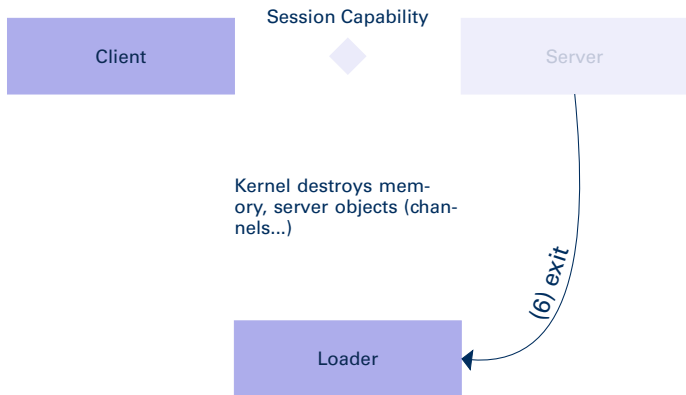
## L4Re: Session Creation



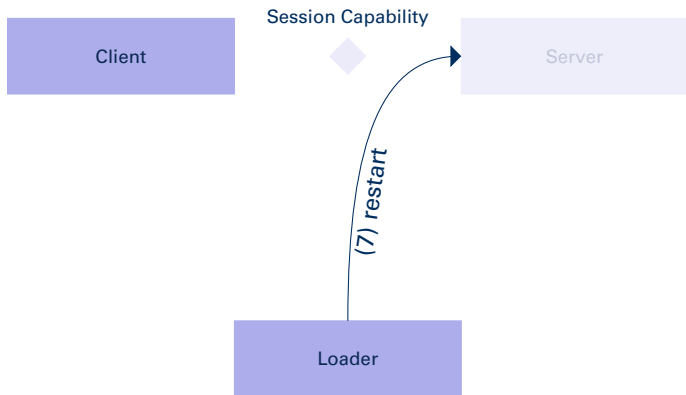
## L4Re: Session Creation



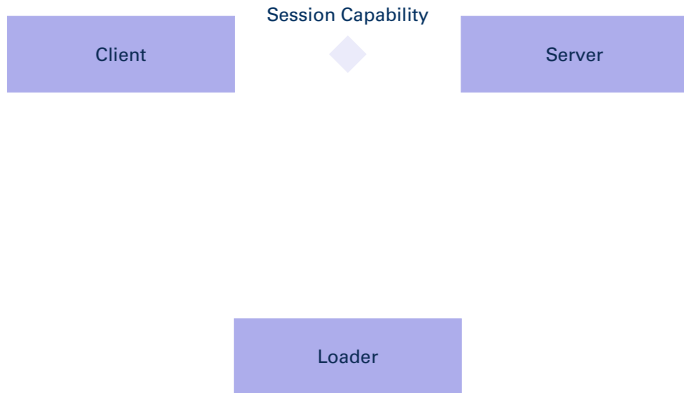
## L4Re: Server Crash



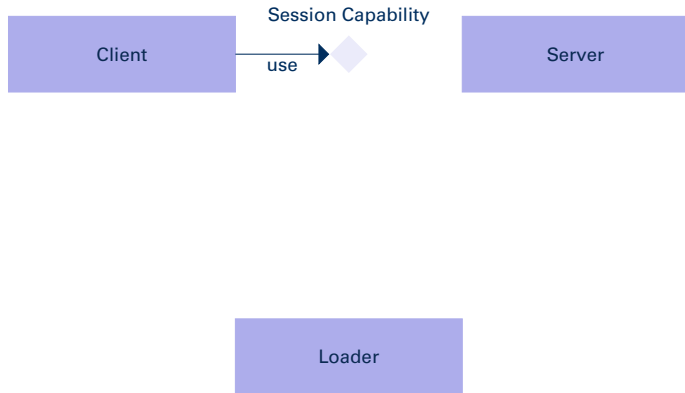
## L4Re: Server Crash



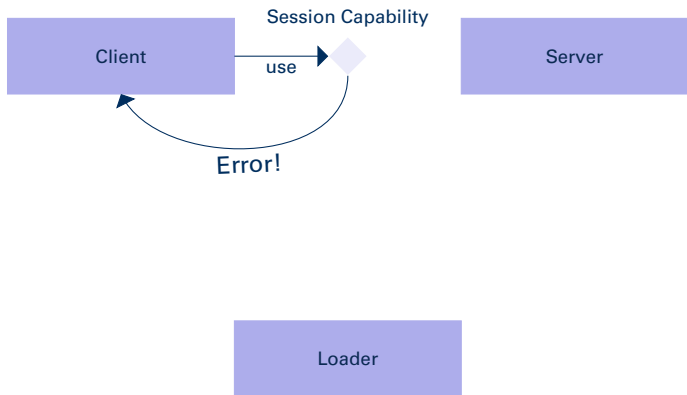
## L4Re: Restarted Server



## L4Re: Restarted Server



## L4Re: Restarted Server



## L4ReAnimator

- Only the application itself can detect that a capability vanished
- Kernel raises *Capability fault*
- Application needs to re-obtain the capability: execute *capability fault handler*
- Capfault handler: application-specific
  - Create new communication channel
  - Restore session state
- Programming model:
  - Capfault handler provided by server implementor
  - Handling transparent for application developer
  - *Semi-transparency*

## L4ReAnimator: Cleanup

- Some channels have resources attached (e.g., frame buffer for graphical console)
- Resource may come from a different resource (e.g., frame buffer from memory manager)
- Resources remain intact (stale) upon crash
- Client ends up using old version of the resource
- Requires additional app-specific knowledge
- *Unmap handler*

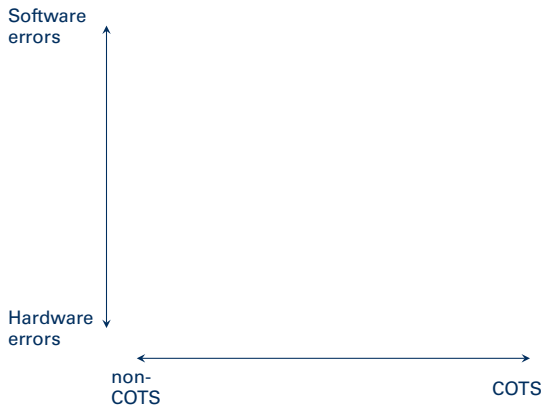
## Summary

- L4ReAnimator
  - Capfault: Clients detect server restarts lazily
  - Capfault Handler: application-specific knowledge on how to regain access to the server
  - Unmap handler: clean up old resources after restart
  
- All these frameworks only deal with software errors.
- What about hardware faults?

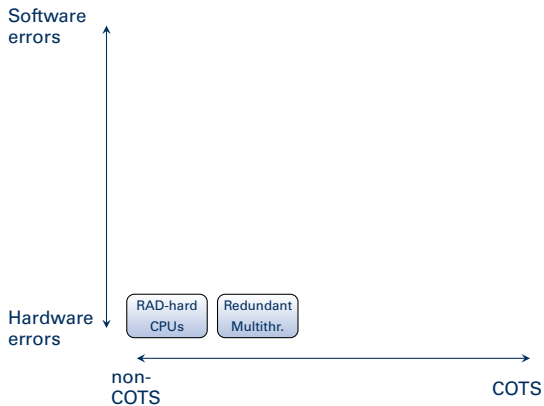
# Hardware Errors

- Hardware fails all the time
- Permanent faults
  - Studies: 2% chance of a DRAM DIMM to fail within a year
  - Hardware ECC cannot catch them all!
  - 5% chance of a disk failure within a year
- Decreasing transistor sizes → increase in rate of transient (soft) errors
  - No longer only a space/aviation problem

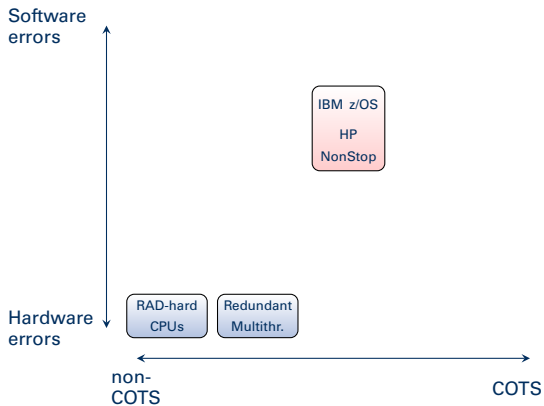
## Fault Tolerance: State of the Union



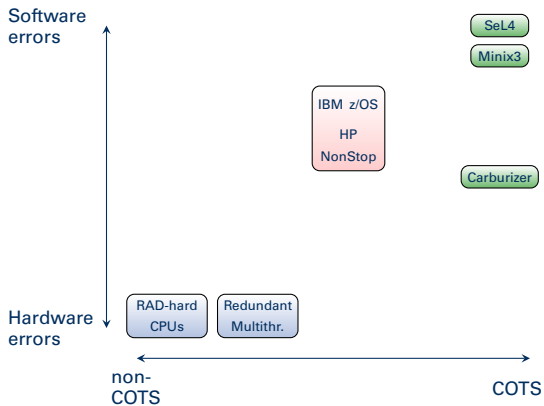
## Fault Tolerance: State of the Union



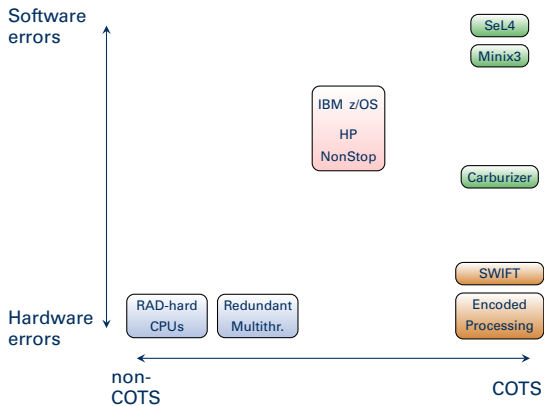
## Fault Tolerance: State of the Union



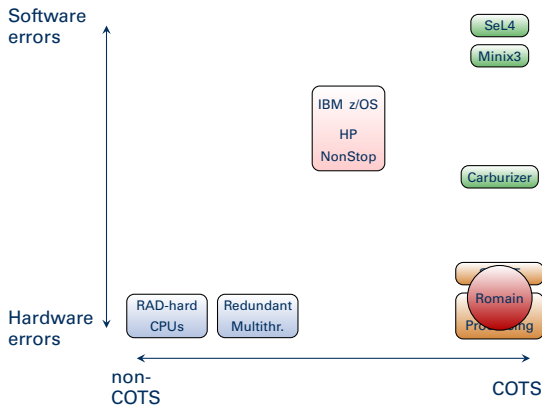
## Fault Tolerance: State of the Union



# Fault Tolerance: State of the Union



# Fault Tolerance: State of the Union



## Process-Level Redundancy [Shye 2007]

### Binary recompilation

- Complex, unprotected compiler
- Architecture-dependent

### System calls for replica synchronization

### Virtual memory fault isolation

- Restricted to Linux user-level programs

## Process-Level Redundancy [Shye 2007]

### Binary recompilation

- Complex, unprotected compiler
- Architecture-dependent

Reuse OS mechanisms

### System calls for replica synchronization

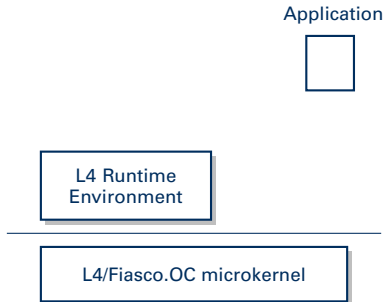
Additional synchronization events

### Virtual memory fault isolation

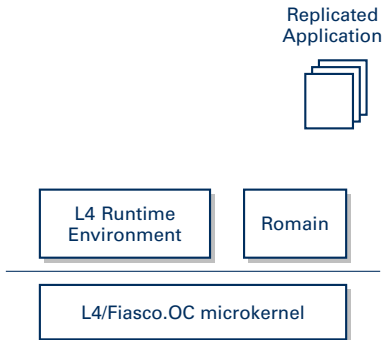
- Restricted to Linux user-level programs

Microkernel-based

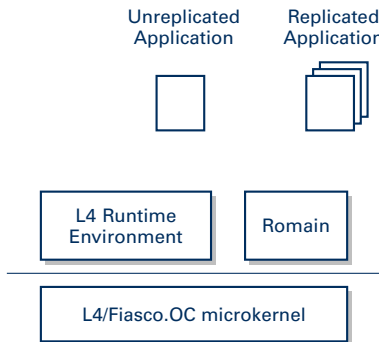
# Transparent Replication as OS Service



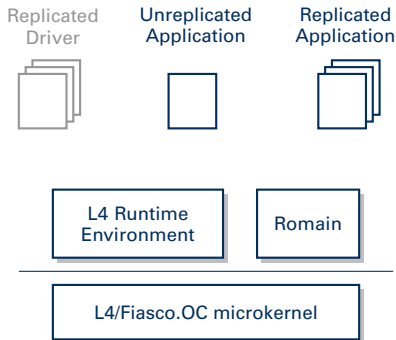
# Transparent Replication as OS Service



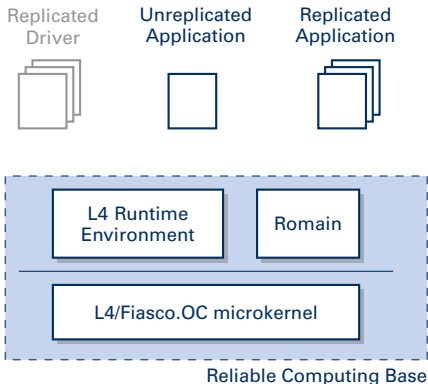
# Transparent Replication as OS Service



# Transparent Replication as OS Service



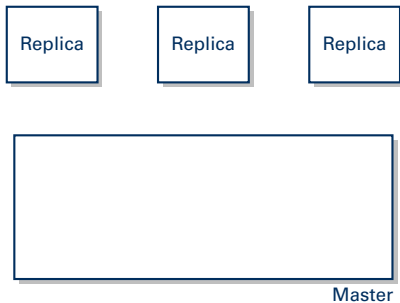
# Transparent Replication as OS Service



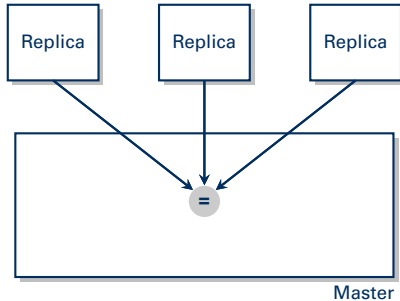
## Romain: Structure



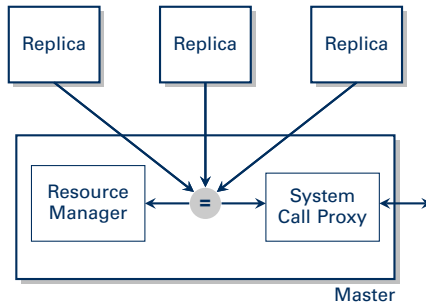
## Romain: Structure



## Romain: Structure



## Romain: Structure



# Resource Management: Capabilities

Replica 1



# Resource Management: Capabilities

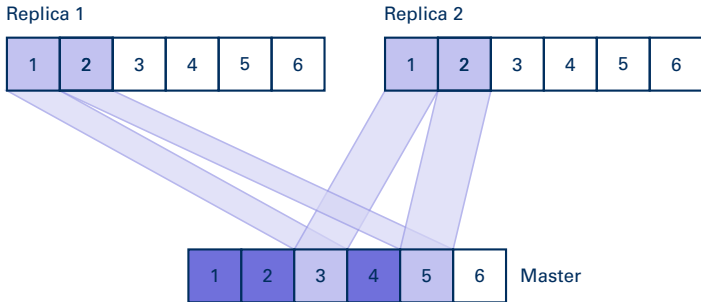
Replica 1



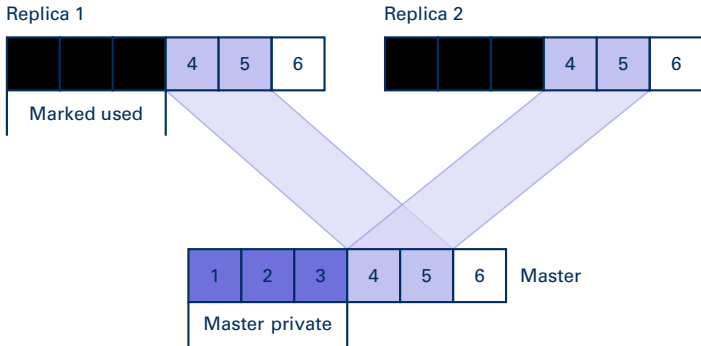
Replica 2



# Resource Management: Capabilities

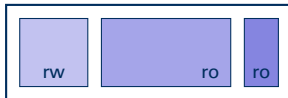


# Partitioned Capability Tables

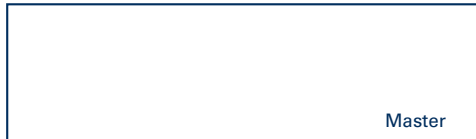


# Replica Memory Management

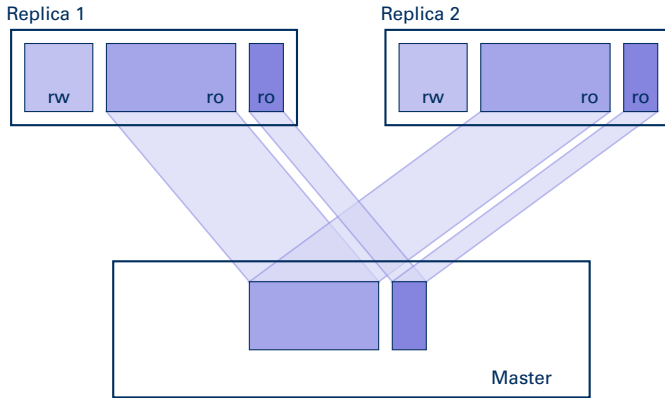
Replica 1



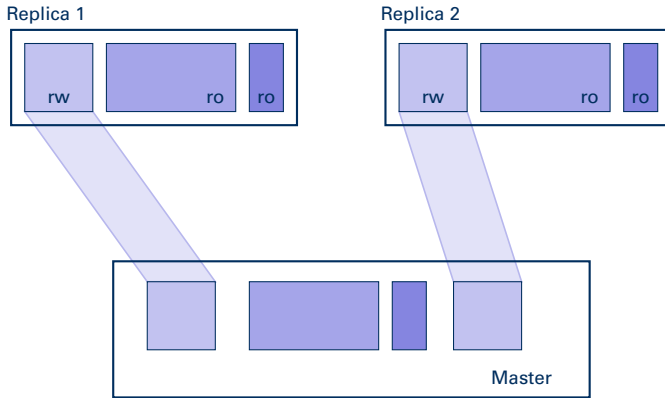
Replica 2



# Replica Memory Management



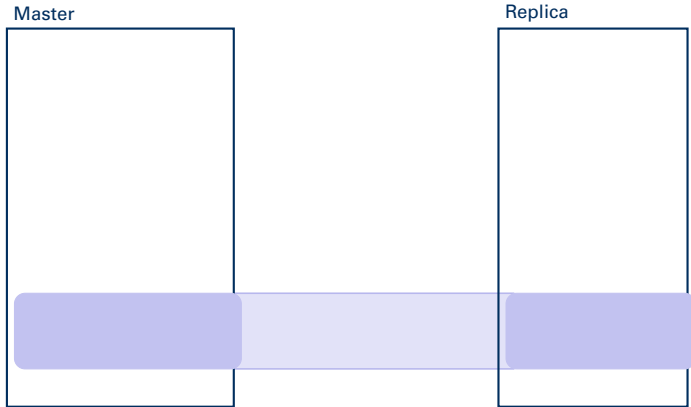
# Replica Memory Management



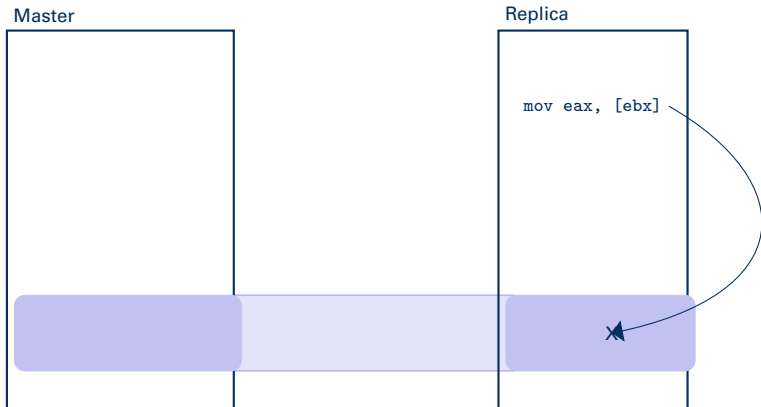
## Shared Memory

- Not in complete control of master
- Standard technique: trap&emulate
  - Execution overhead (x100 - x1000)
  - Adds complexity to RCB
    - Disassembler 6,000 LoC
    - Tiny emulator 500 LoC
- Our implementation: copy & execute

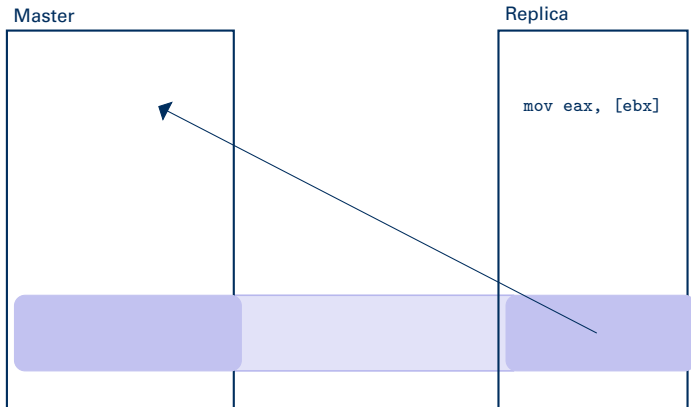
# Copy&Execute



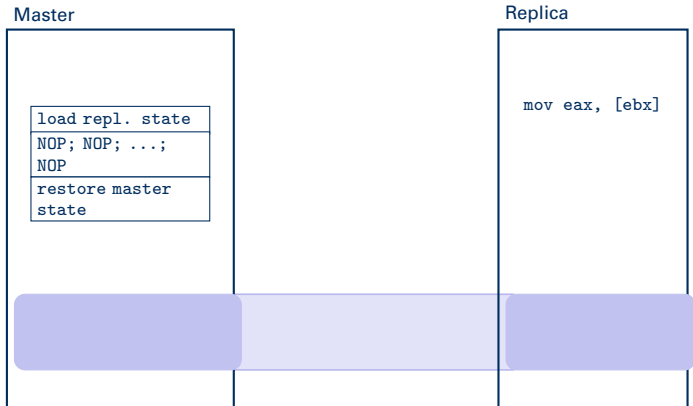
# Copy&Execute



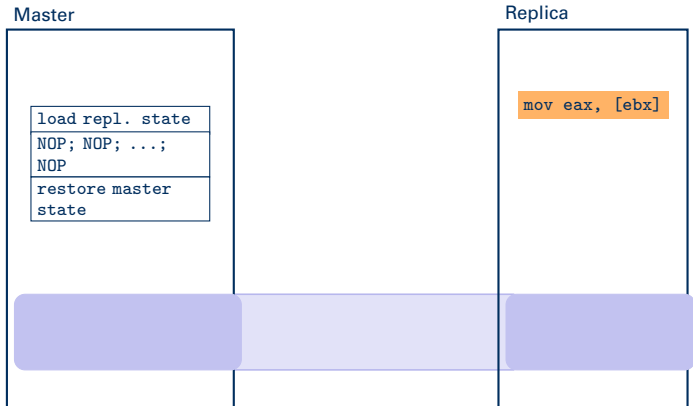
# Copy&Execute



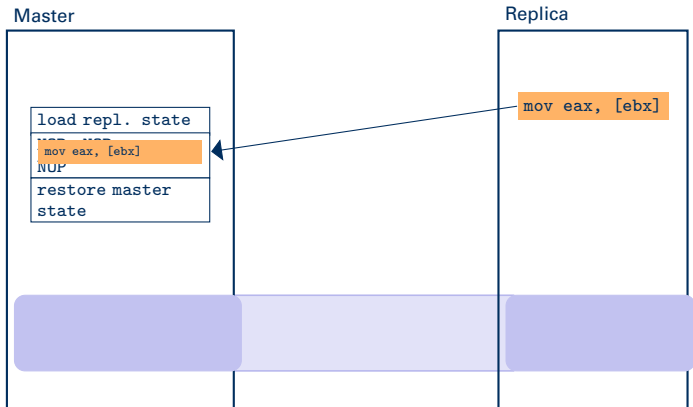
# Copy&Execute



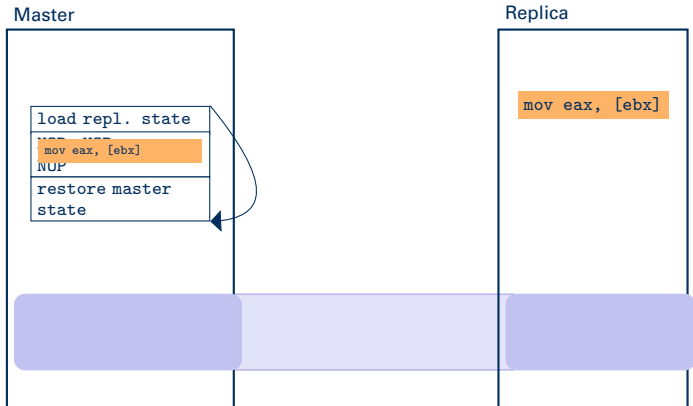
# Copy&Execute



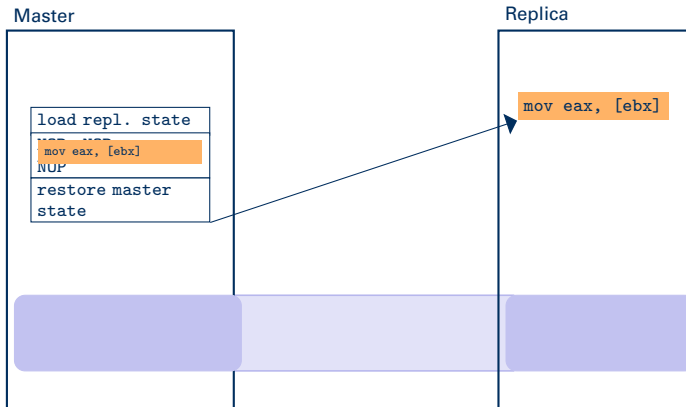
# Copy&Execute



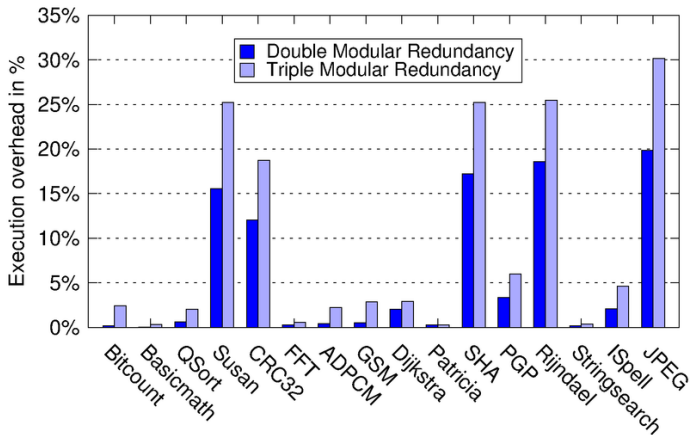
# Copy&Execute



# Copy&Execute



## Overhead vs. Unreplicated Execution



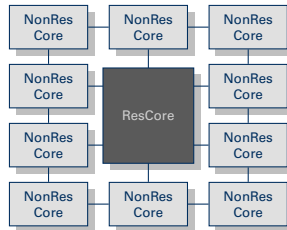
## Romain Lines of Code

Base code (main, logging, locking)	325
Application loader	375
Replica manager	628
Redundancy	153
Memory manager	445
System call proxy	311
Shared memory	281
<b>Total</b>	<b>2,518</b>
Fault injector	668
GDB server stub	1,304

## Hardening the RCB

- **We need:** Dedicated mechanisms to protect the RCB (HW or SW)
- **We have:** Full control over software
- Use FT-encoding compiler?
  - Has not been done for kernel code yet
- RAD-hardened hardware?
  - Too expensive

Why not split cores into resilient and non-resilient ones?



## Summary

- OS-level techniques to tolerate SW and HW faults
- Address-space isolation
- Microreboots
- Various ways of handling session state
- Replication against hardware errors

## Further Reading

- **Minix3:** Jorrit Herder, Ben Gras,, Philip Homburg, Andrew S. Tanenbaum: *Fault Isolation for Device Drivers*, DSN 2009
- **CuriOS:** Francis M. David, Ellick M. Chan, Jeffrey C. Carlyle and Roy H. Campbell *CuriOS: Improving Reliability through Operating System Structure*, OSDI 2008
- **L4ReAnimator:** Dirk Vogt, Björn Döbel, Adam Lackorzynski: *Stay strong, stay safe: Enhancing Reliability of a Secure Operating System*, IIDS 2010
- **PLR:** Alex Shye, Tipp Moseley, Vijay Janapa Reddi, Joseh Blomsted, Ramesh Peri: *Using Process-Level Redundancy to Exploit Multiple Cores for Transient Fault Tolerance*, DSN 2007
- **Romain:**
  - Björn Döbel, Hermann Härtig, Michael Engel: *Operating System Support for Redundant Multithreading*, EMSOFT 2012
  - Björn Döbel, Hermann Härtig: *Who watches the watchmen? – Protecting Operating System Reliability Mechanisms*, HotDep 2012