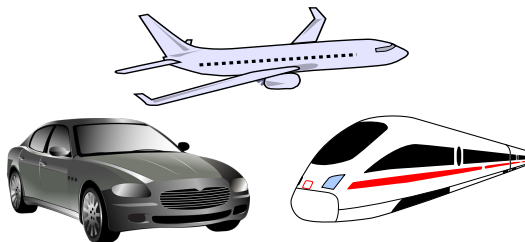


# The Hierarchical Microkernel: A Flexible and Robust OS Architecture

Stefan Winter, Martin Tsarev, Neeraj Suri  
DEEDS Group, TU Darmstadt  
moduli-os@deeds.informatik.tu-darmstadt.de

8th November 2013

# My PhD research focus



## Robustness

Correct operation despite

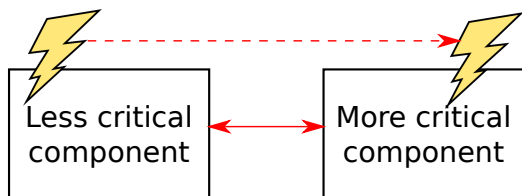
- invalid inputs
- stressful environmental conditions

Problem  
Statement

Existing  
Approaches

Proposed Solution

Summary



## Robustness

Correct operation despite

- invalid inputs
- stressful environmental conditions

Problem  
Statement

Existing  
Approaches

Proposed Solution

Summary

# How (most) software systems evolve

## Composition

- static linking
- dynamic linking

## Problems

- static: system down time
- dynamic...?

# How (most) software systems evolve

## Composition

- static linking
- dynamic linking

## Problems

- static: system down time
- dynamic...?

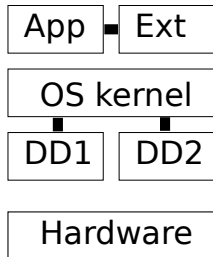
# Monolithic compositions do (not) support robust evolution

Problem  
Statement

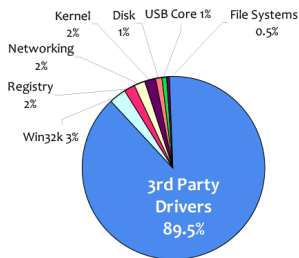
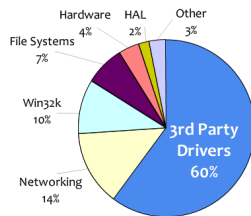
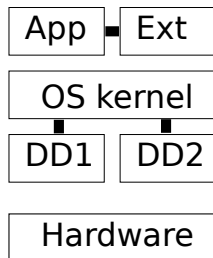
Existing  
Approaches

Proposed Solution

Summary



# Monolithic compositions do (not) support robust evolution



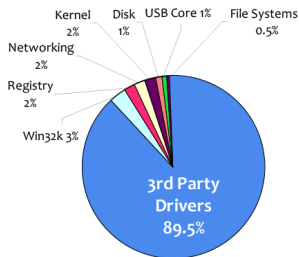
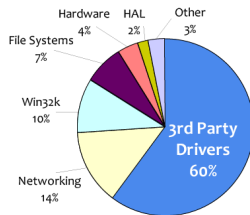
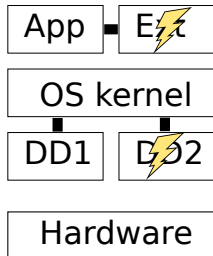
Problem  
Statement

Existing  
Approaches

Proposed Solution

Summary

# Monolithic compositions do (not) support robust evolution



Problem  
Statement

Existing  
Approaches

Proposed Solution

Summary



# Monolithic compositions do (not) support robust evolution

Problem  
Statement

Existing  
Approaches

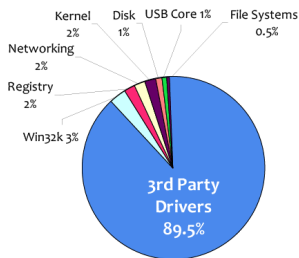
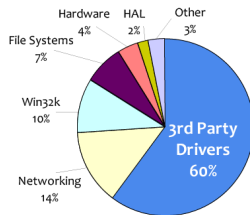
Proposed Solution

Summary

App	Ext
-----	-----

OS kernel	
DD1	DD2

Hardware
----------



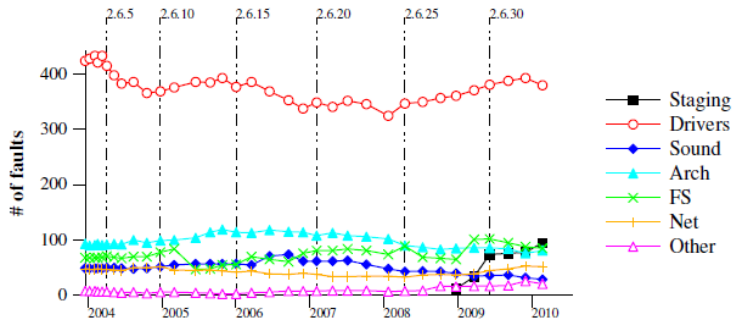
# Extensions are difficult to get right

Problem  
Statement

Existing  
Approaches

Proposed Solution

Summary



© Palix et al., *Faults in Linux: ten years later*, ASPLOS'11

- ① components of different criticality
- ② components of different origin (COTS/SOUP)
- ③ uniform privilege
- ④ complex and volatile interfaces evolve independently

# Solution 1: Sandboxing

- Nooks: Linux driver sandboxing<sup>1</sup>
- Microdrivers<sup>2</sup>: performance-critical code in the kernel
- BGI: Byte-Granularity Isolation<sup>3</sup>
- (L)XFI: Windows/Linux in-kernel fault isolation<sup>4,5</sup>

## Sandboxing issues

- available?
- working?
- *co-evolution* with OS required

---

<sup>1</sup>Swift et al.: *Improving the Reliability of Commodity Operating Systems*, SOSP'03

<sup>2</sup>Ganapathy et al.: *The design and implementation of microdrivers*, ASPLOS'08

<sup>3</sup>Castro et al.: *Fast Byte-Granularity Software Fault Isolation*, SOSP'09

<sup>4</sup>Erlingsson et al.: *XFI: software guards for system address spaces*, OSDI'06

<sup>5</sup>Mao et al.: *Software fault isolation with API integrity and multi-principal modules*, SOSP'11

# Solution 1: Sandboxing

- Nooks: Linux driver sandboxing<sup>1</sup>
- Microdrivers<sup>2</sup>: performance-critical code in the kernel
- BGI: Byte-Granularity Isolation<sup>3</sup>
- (L)XFI: Windows/Linux in-kernel fault isolation<sup>4,5</sup>

## Sandboxing issues

- available?
- working?
- *co-evolution* with OS required

---

<sup>1</sup>Swift et al.: *Improving the Reliability of Commodity Operating Systems*, SOSP'03

<sup>2</sup>Ganapathy et al.: *The design and implementation of microdrivers*, ASPLOS'08

<sup>3</sup>Castro et al.: *Fast Byte-Granularity Software Fault Isolation*, SOSP'09

<sup>4</sup>Erlingsson et al.: *XFI: software guards for system address spaces*, OSDI'06

<sup>5</sup>Mao et al.: *Software fault isolation with API integrity and multi-principal modules*, SOSP'11

# Solution 2: Isolation by design

- Virtual Machines<sup>6</sup>: high redundancy
- Singularity<sup>7</sup>: type safety, limited runtime protection
- Microkernels<sup>8</sup>

---

<sup>6</sup>LeVasseur et al.: *Unmodified device driver reuse and improved system dependability via virtual machines*, OSDI'04

<sup>7</sup>Hunt et al.: *Broad New OS Research: Challenges and Opportunities*, HotOS'05

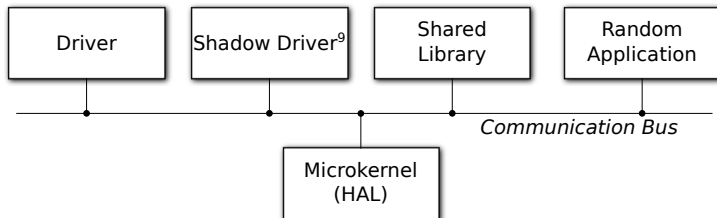
<sup>8</sup>Herder et al.: *Fault isolation for device drivers*, DSN'09

# Proposed solution: “Hierarchical” $\mu$ -kernel

Two core concepts:

- 1 broadcast IPC
- 2 recursive system (de)composition

# Broadcast communication and scalability



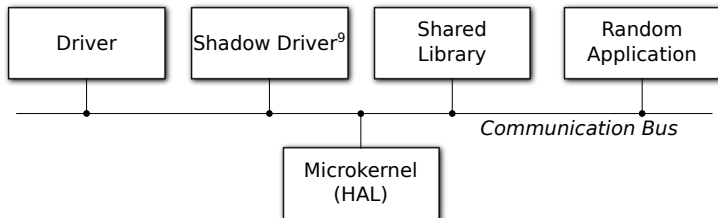
## HM building blocks

- **Modules:** small isolated executable entities
- **Buses:** broadcast message-passing for inter-module communication (in software)

<sup>9</sup>Swift et al.: Recovering device drivers, ACM TOCS 24 4/2006



# Broadcast communication and scalability



## Pros

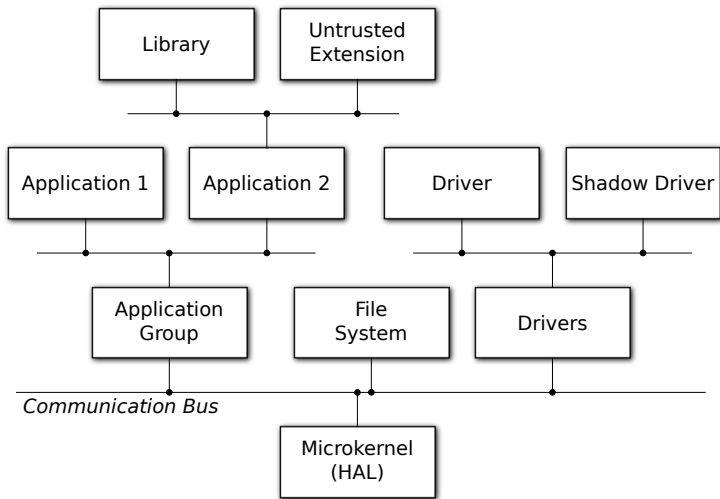
- safe evolution
- no interposition mechanisms required for reconfiguration
- unbounded scalability?

## Cons

- Congestion
- Bus as SPOF
- Confidentiality?  
Availability?

<sup>9</sup>Swift et al.: Recovering device drivers, ACM TOCS 24 4/2006

# Mitigating the downsides of broadcast: Hierarchical (de)composition



# Hierarchical (de)composition

## Hierarchy

**Parent/Child relation** across modules and buses:

- Manage children
  - multiplex resources provided by lower layer
  - provide “system calls”

→ Trust parents

## Pros

- Broadcast scope restriction
- Management load distribution
- Distance from kernel reflects degree of distrust

## Cons

- Communication overheads (routing)
- Hierarchy emulation on binary privilege architectures

## Problems

- ① components of different criticality
- ② components of different origin (COTS/SOUP)
- ③ uniform privilege
- ④ complex and volatile interfaces evolve independently

## Proposal

- *localized* broadcast communication
- *fair* management overhead distribution
- more fine-grained trust/overhead trade-off