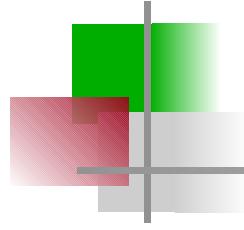


# NIZZA

Hermann Härtig  
TU Dresden ♦OS  
August 2004

Dresden





# dresden – past, present and future

capital of saxony

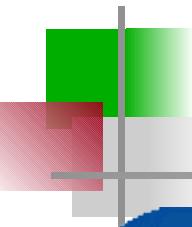
more than 800 years history

half million inhabitants

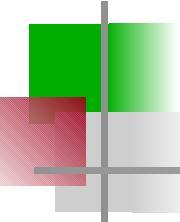
center of saxon silicon valley



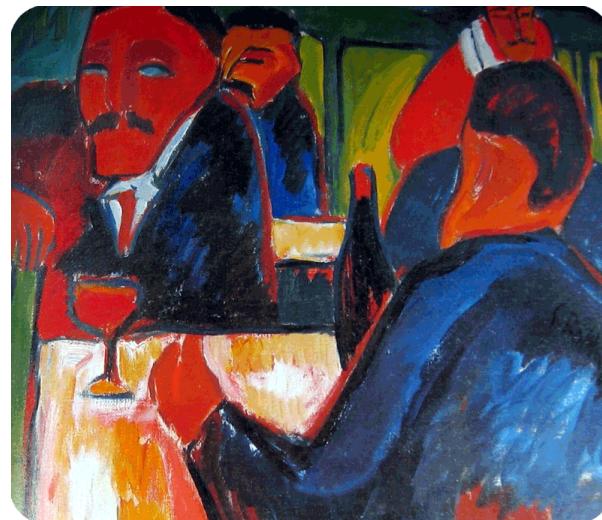
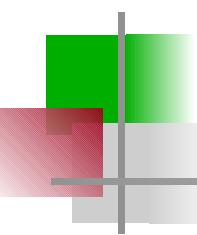
# dresden - impressions



# dresden – art collections

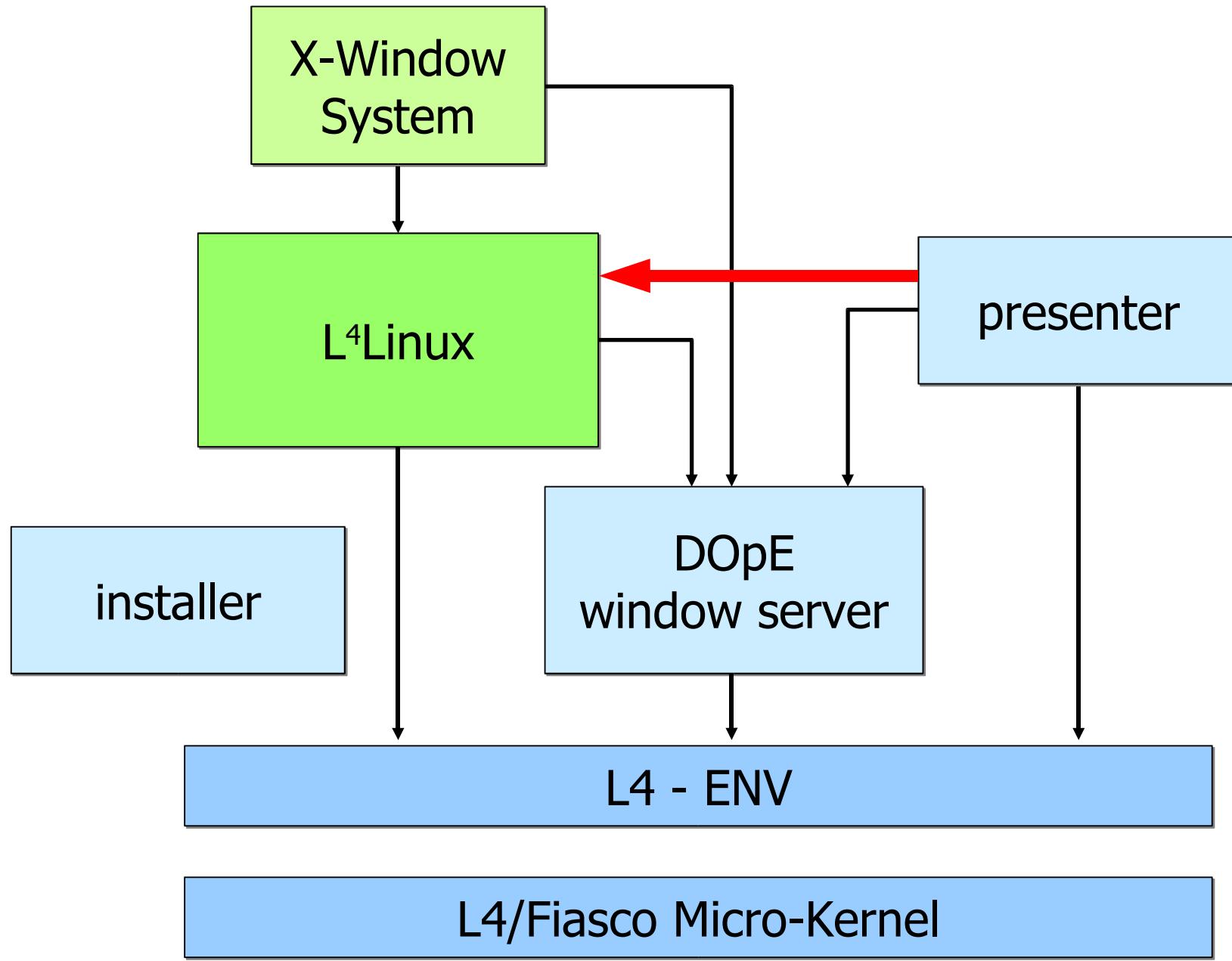


# dresden – home of expressionism

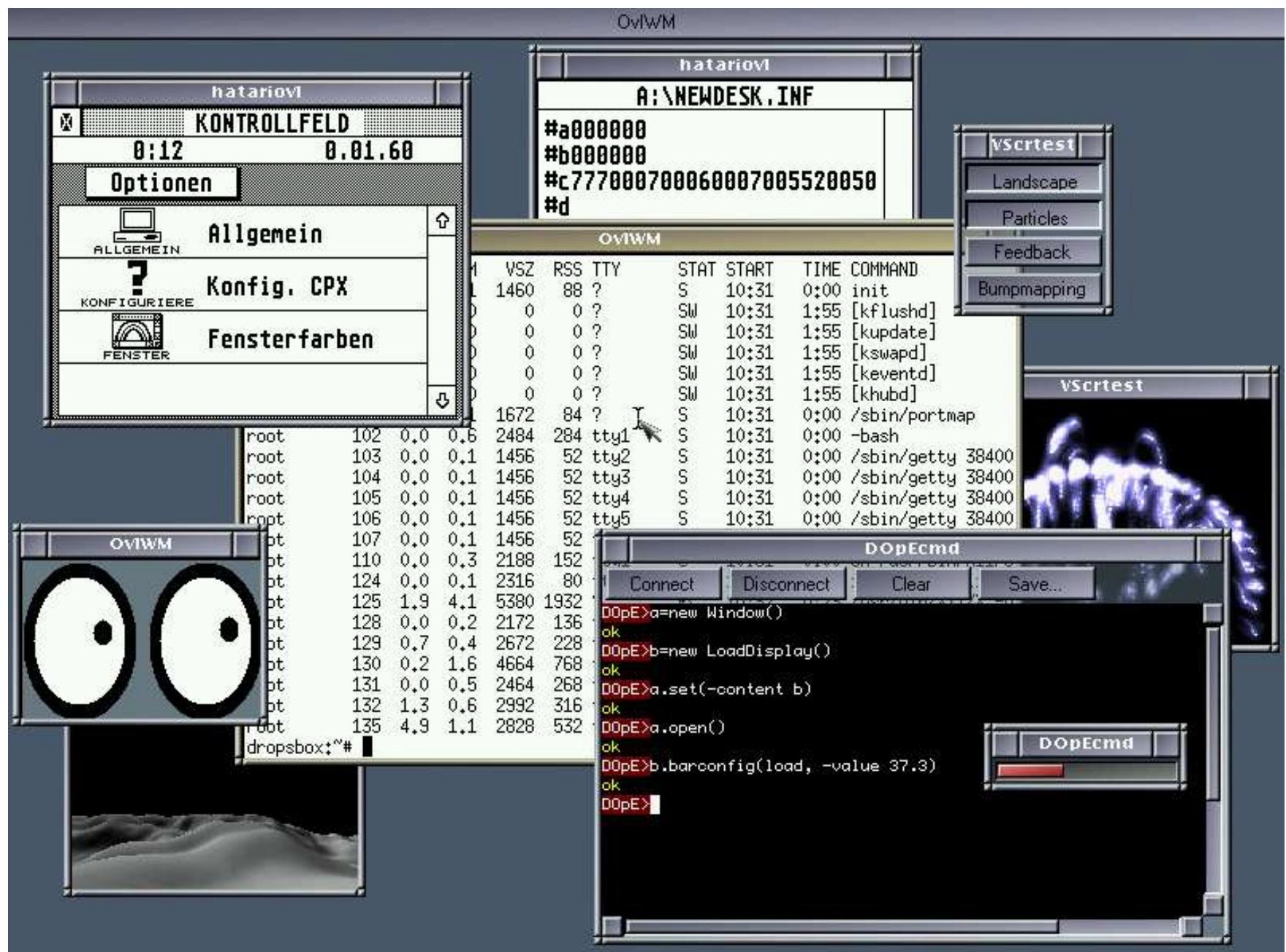


# inventions - made in dresden

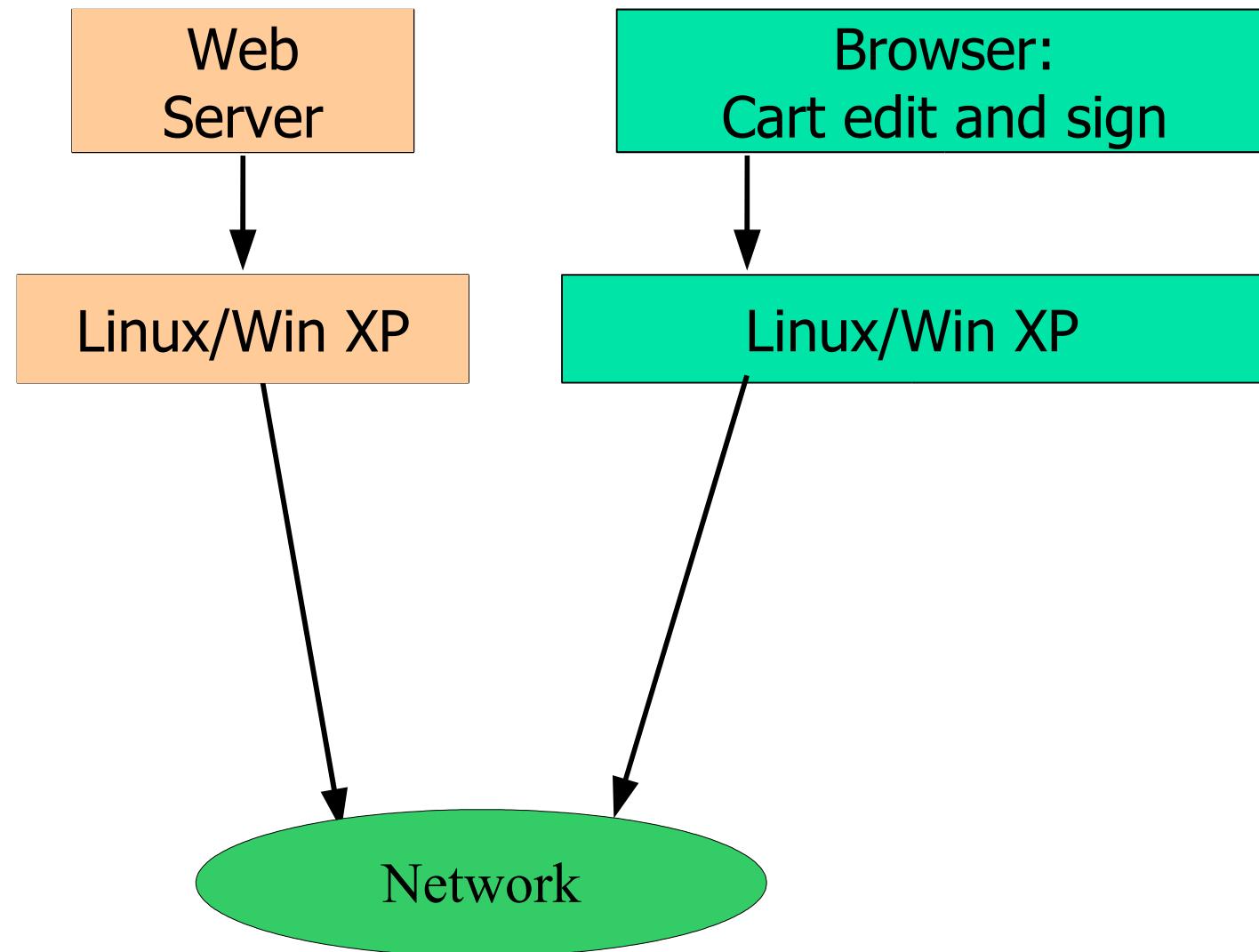


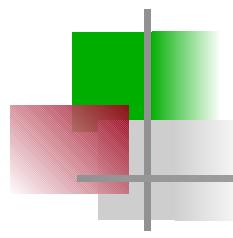


# Screen Shot

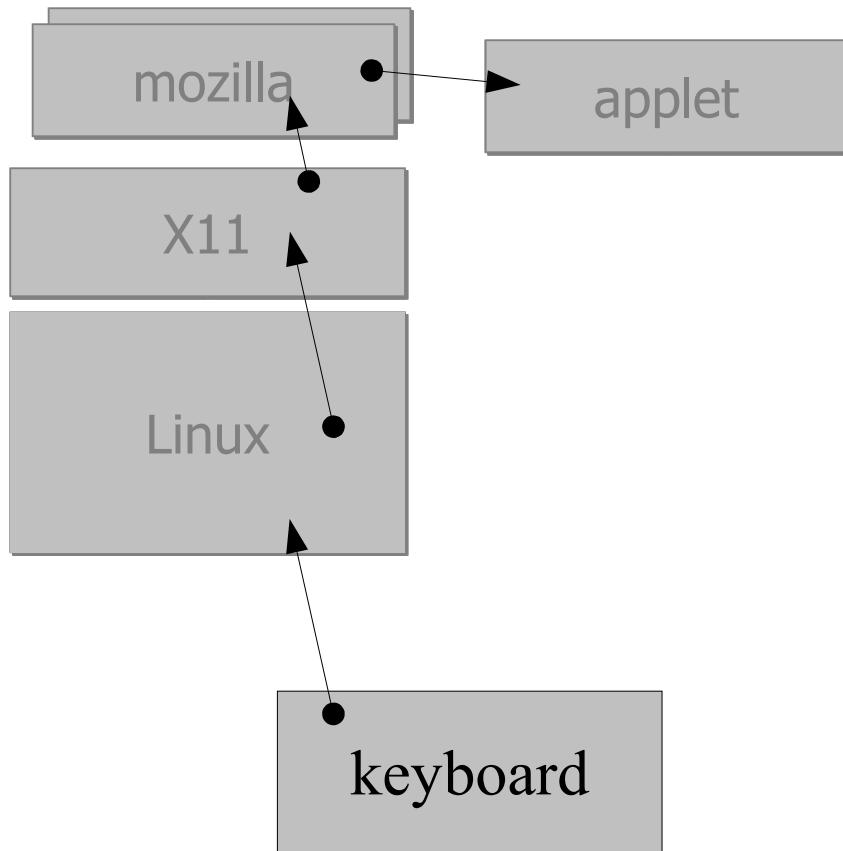


# Internet Transaction





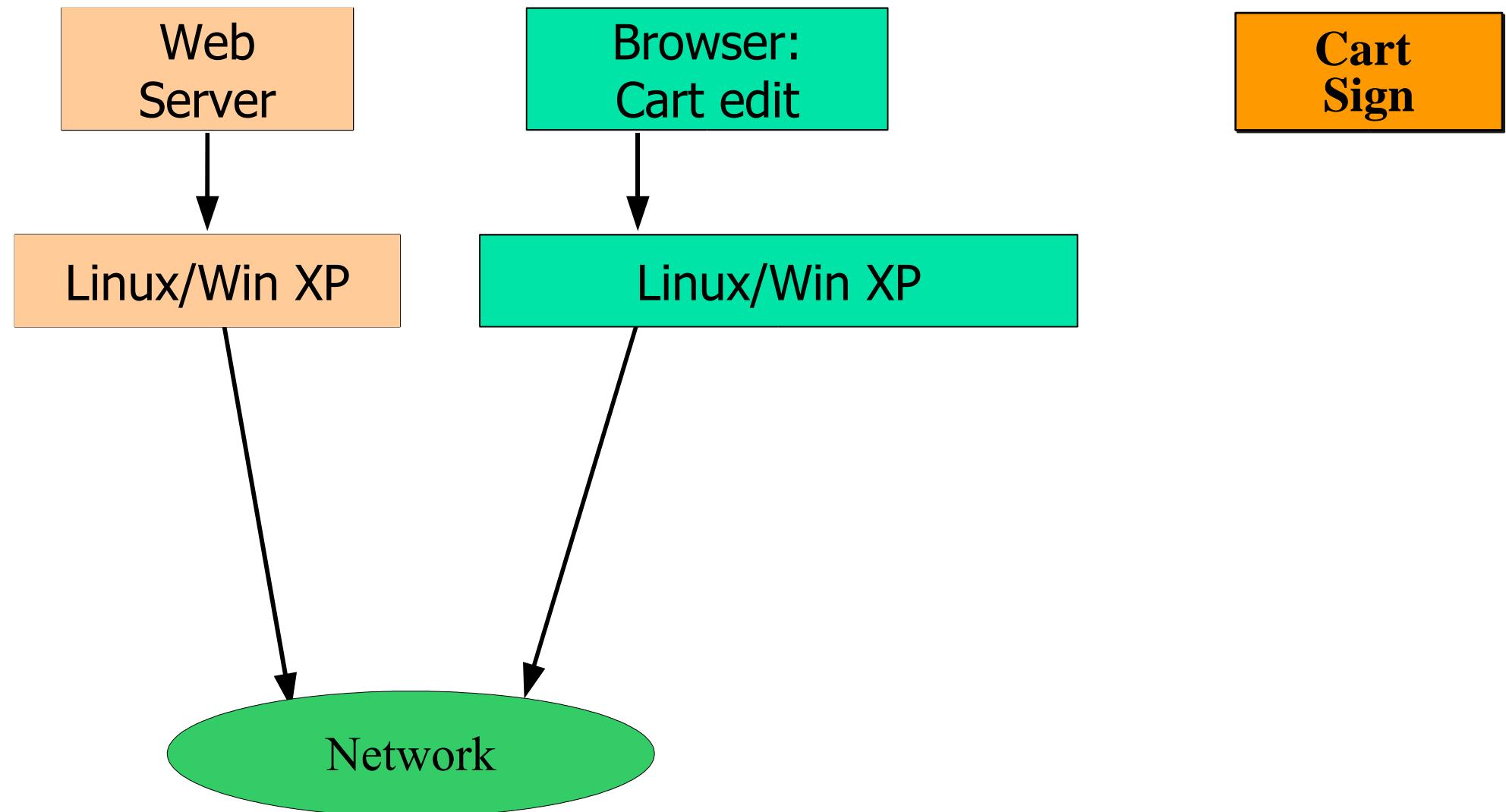
# Your password(s), credit card number, ...



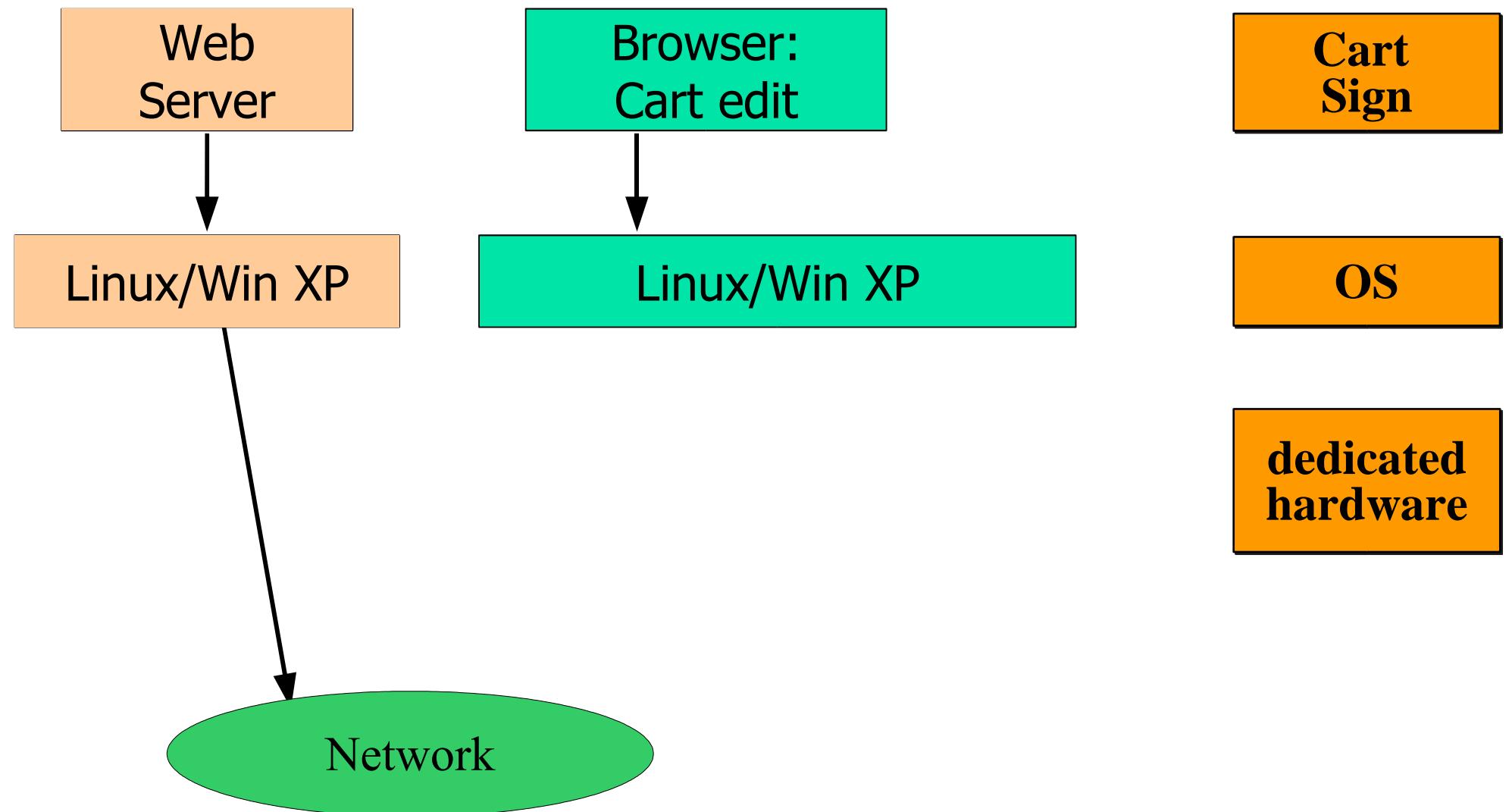
see:

Understanding Data Lifetime  
via Whole System Simulation  
Jim Chow, Ben Pfaff, Tal  
Garfinkel, Kevin Christopher,  
and Mendel Rosenblum,  
Stanford University  
Usenix Security 04

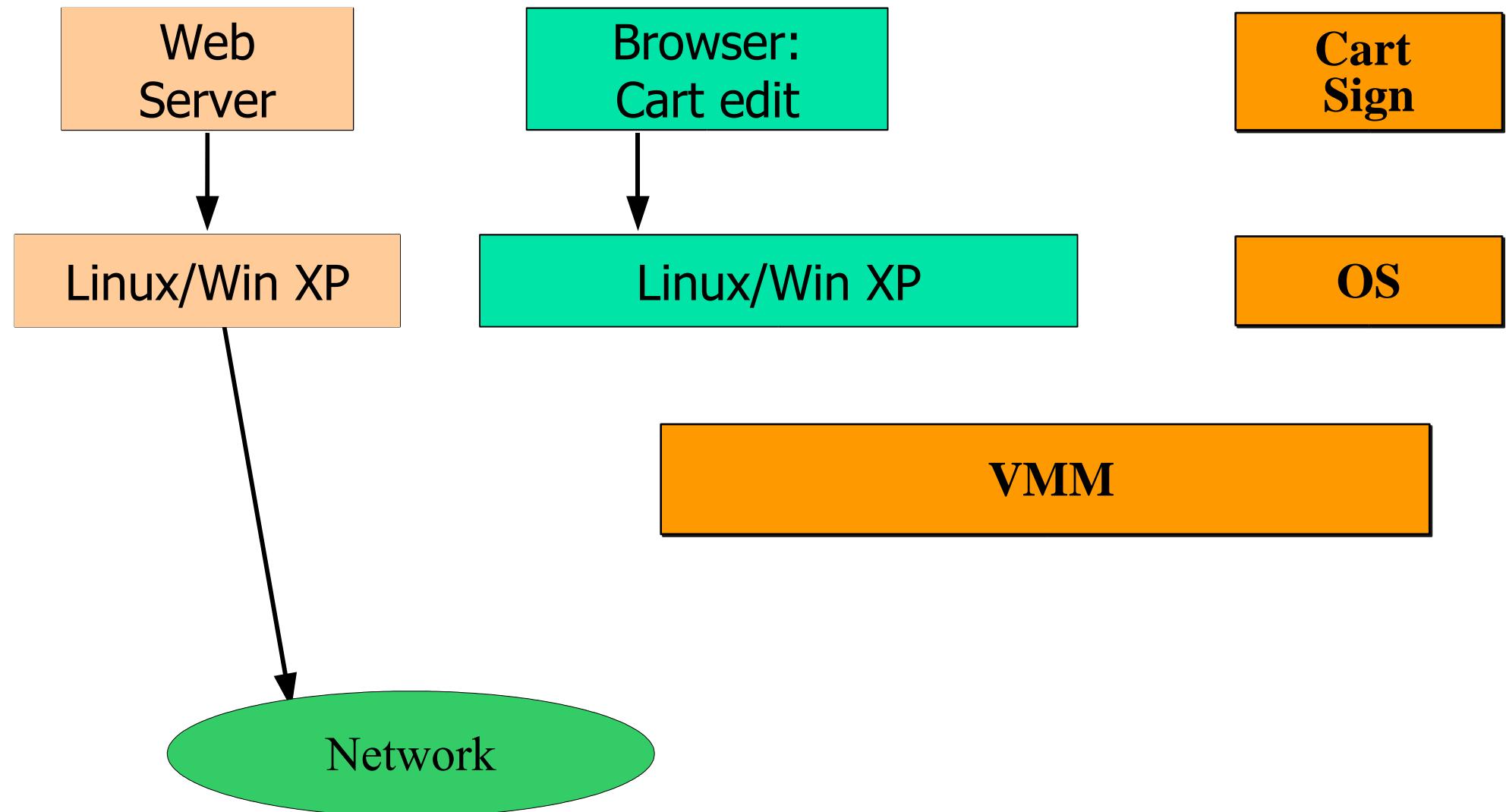
# Lenin's Transaction



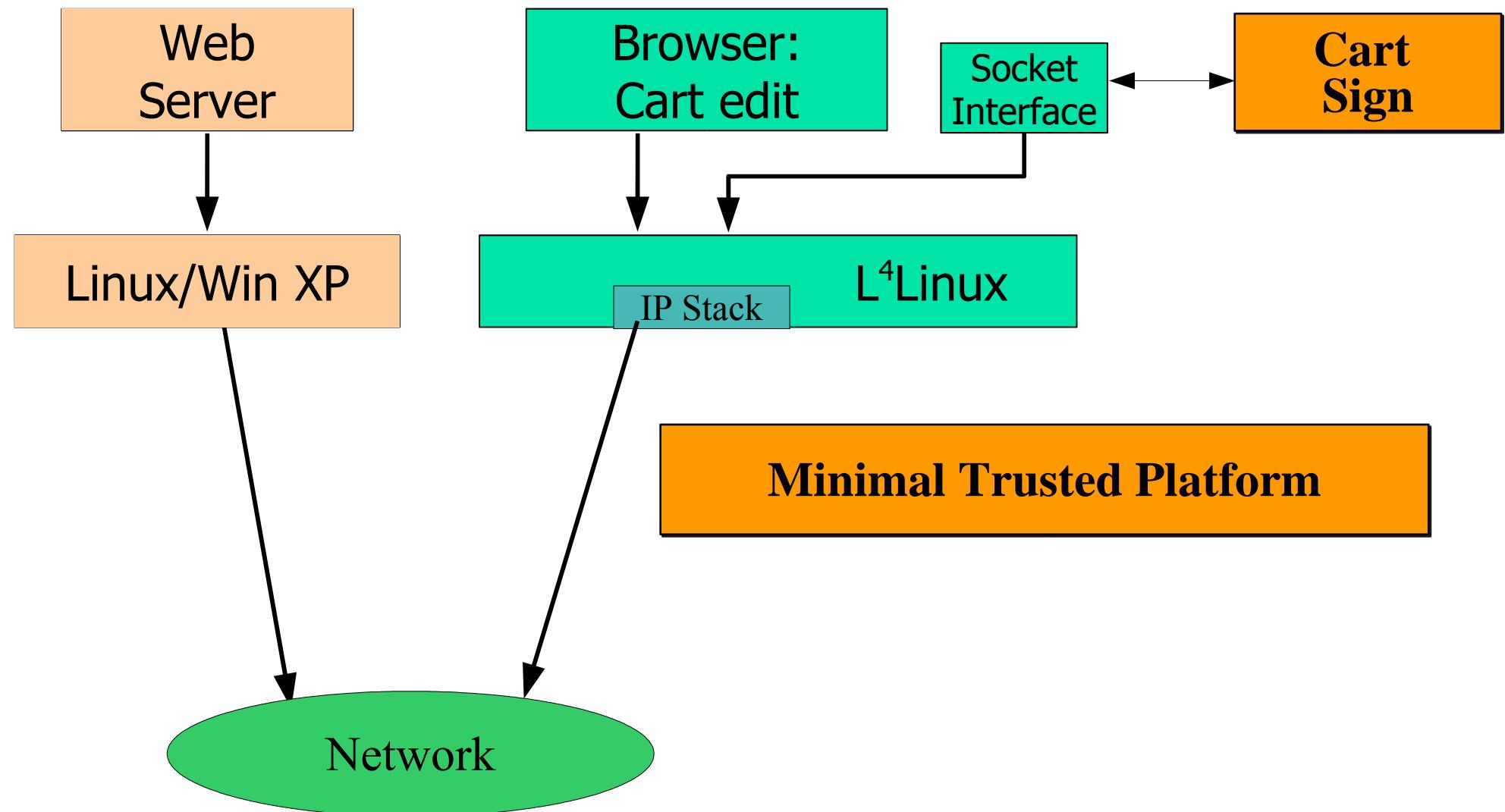
# Lenin's Transaction



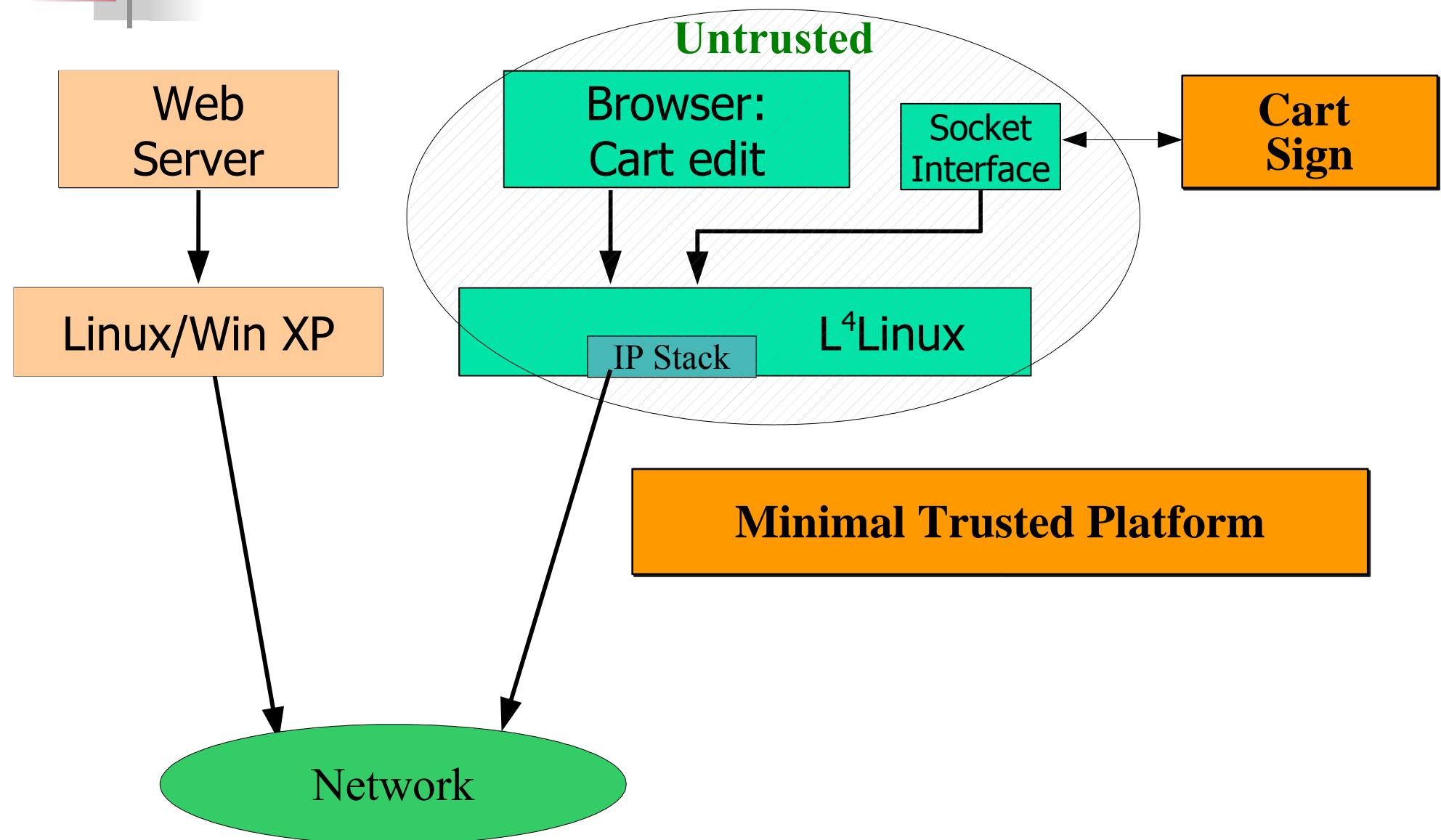
# Lenin's Transaction



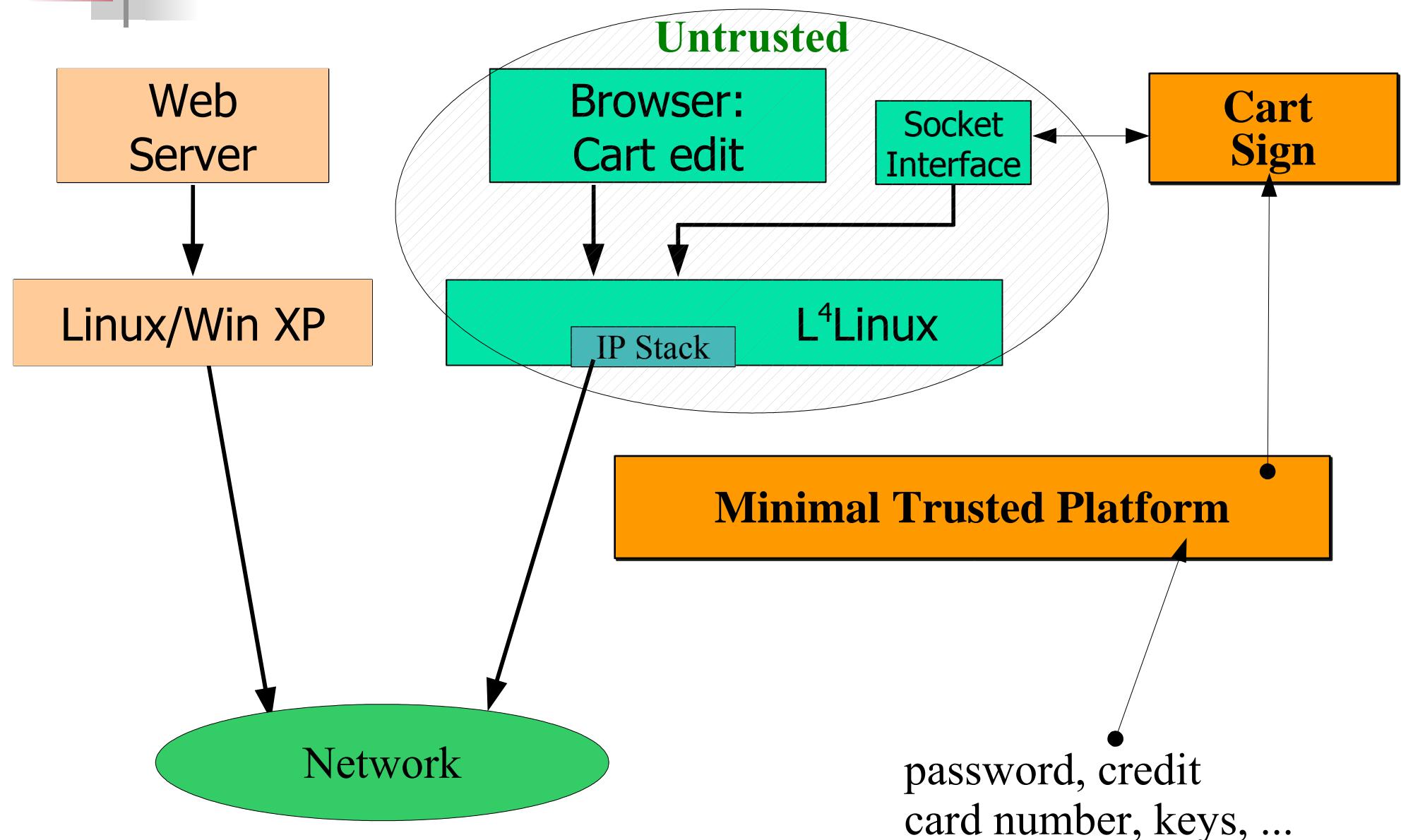
# Lenin's Transaction



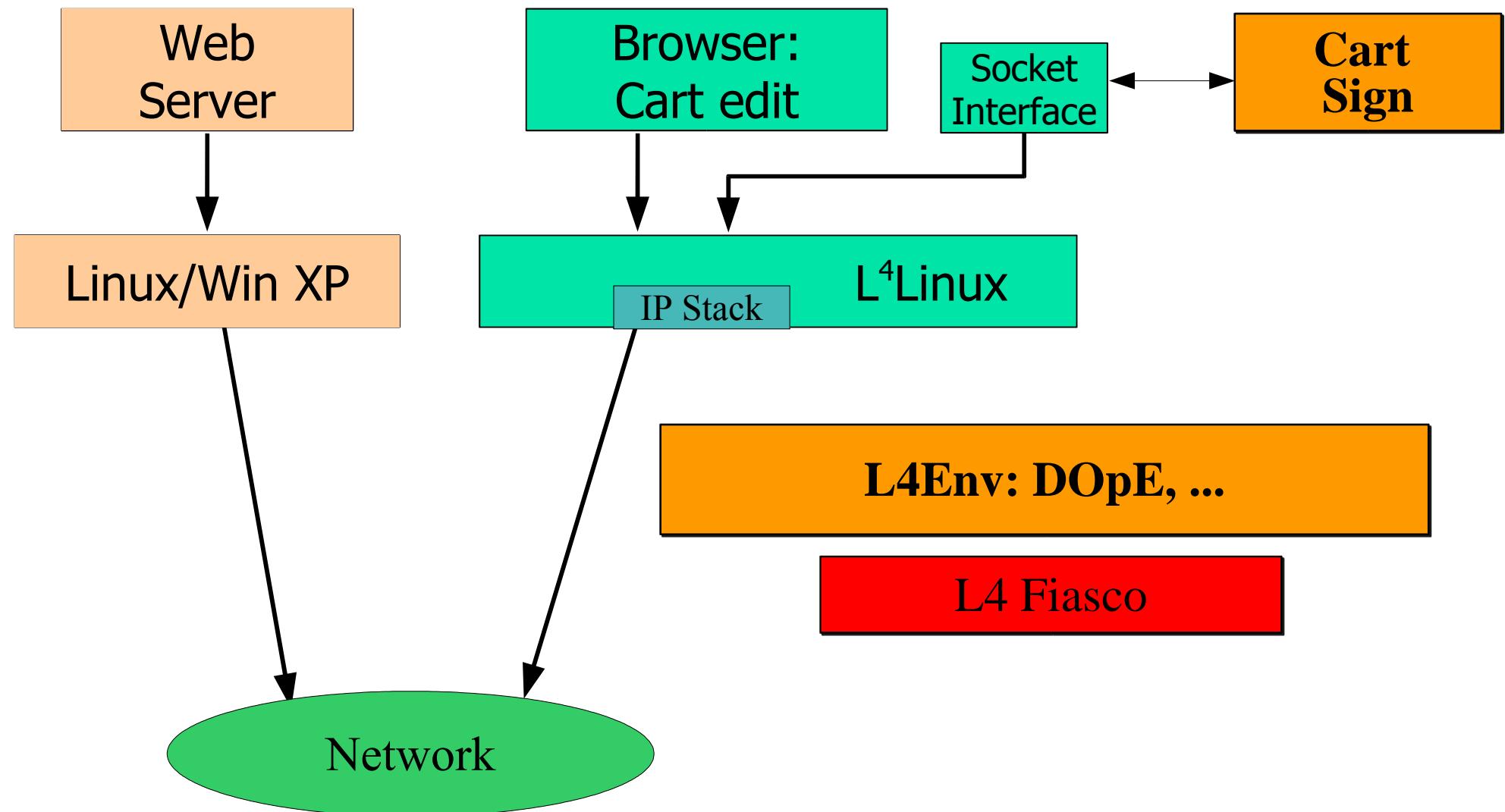
# Lenin's Transaction



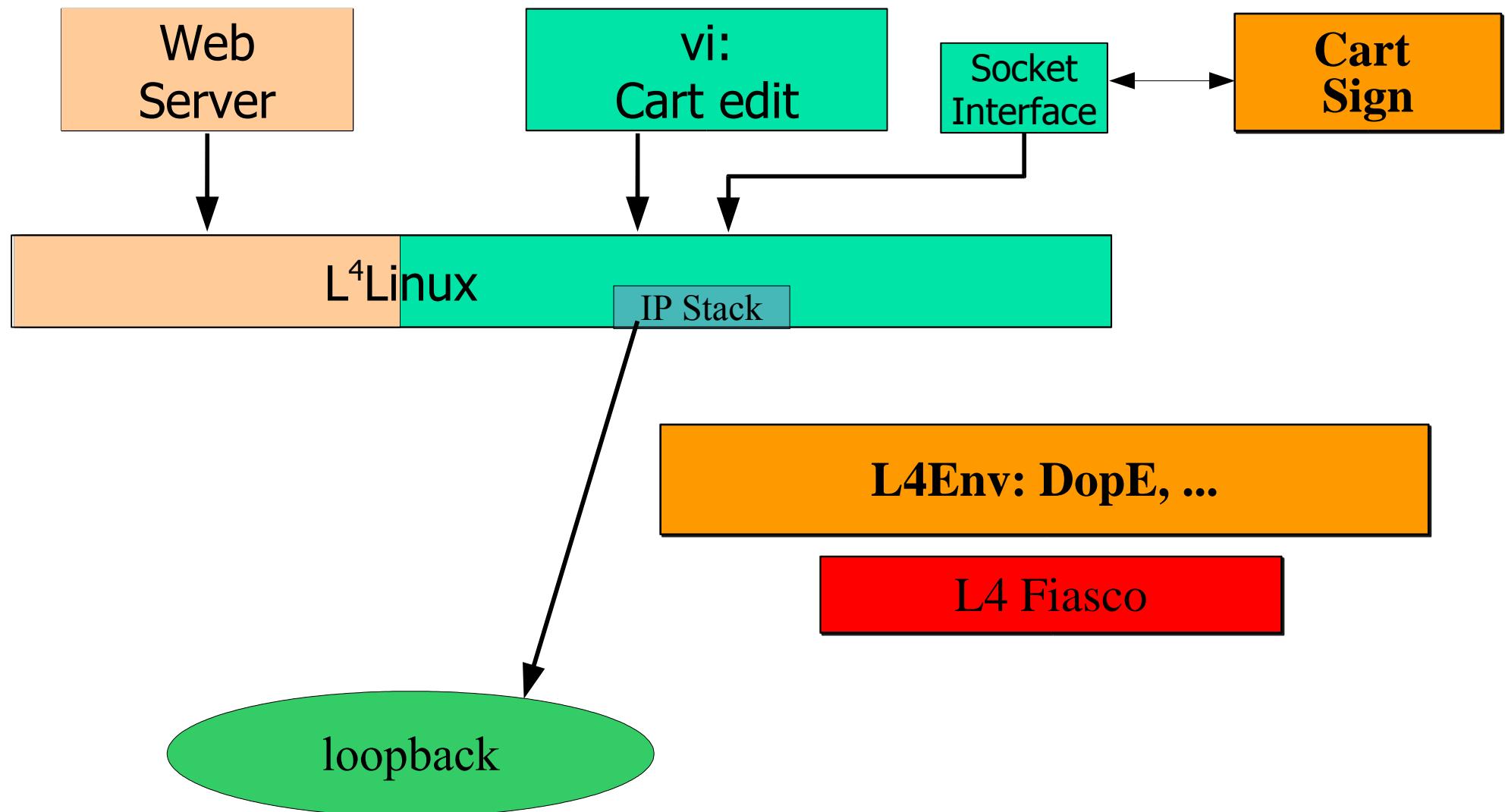
# Lenin's Transaction



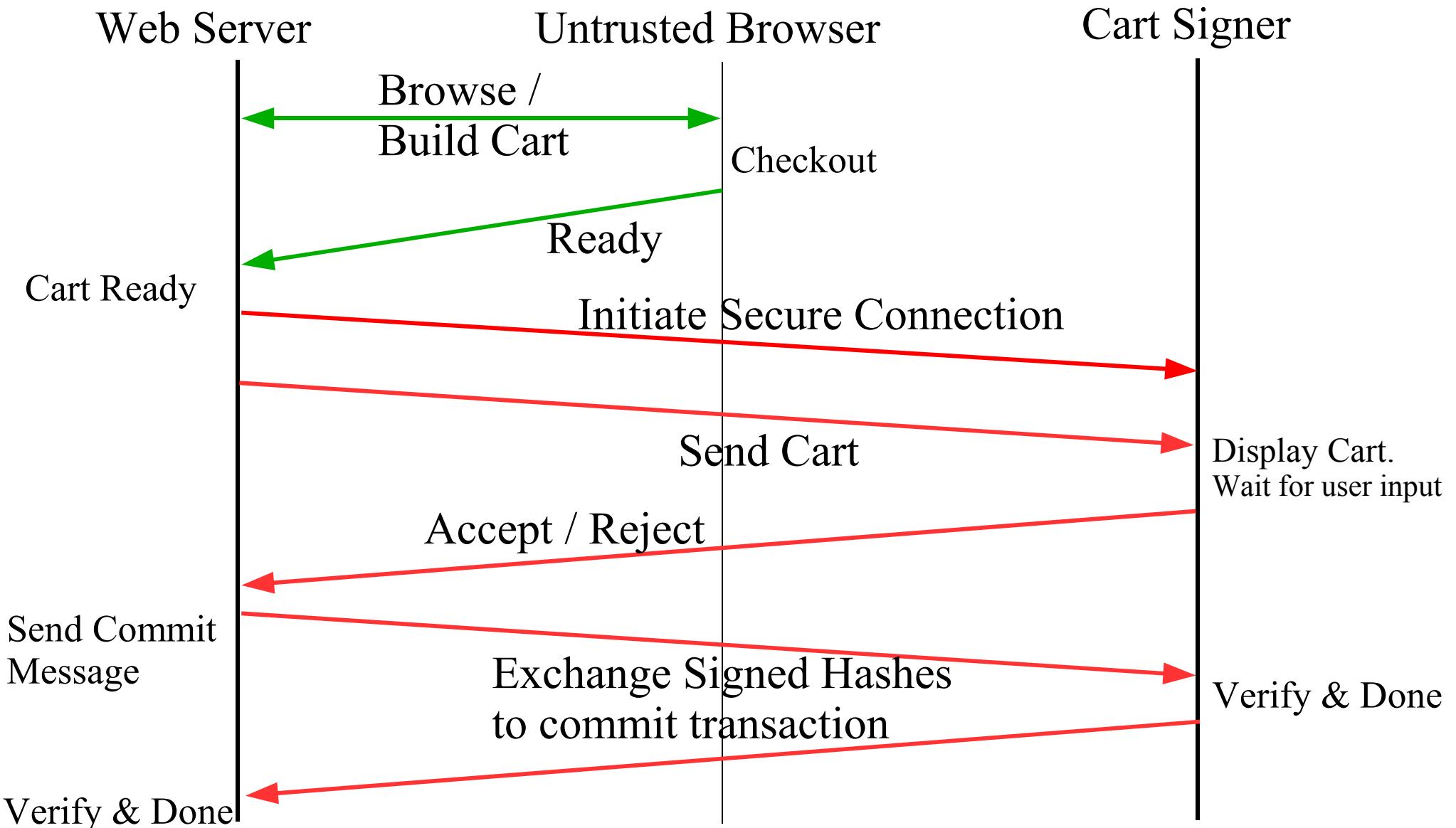
# Lenin's Demo



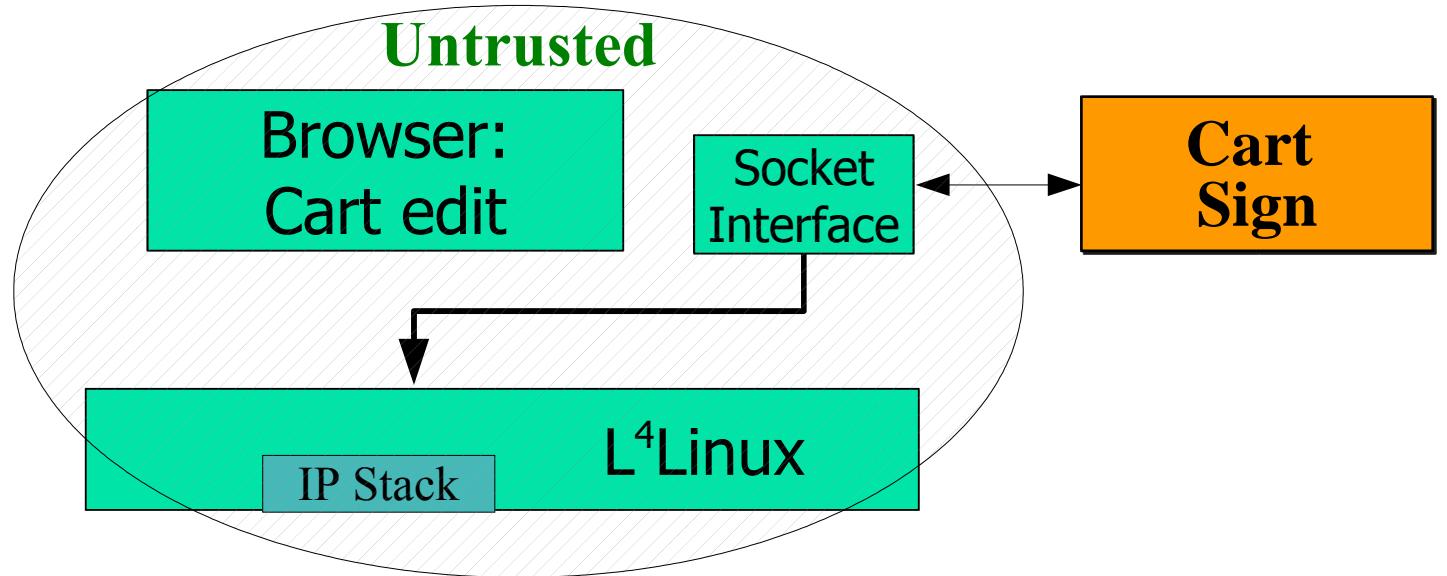
# Lenin's Demo



# Message Sequence

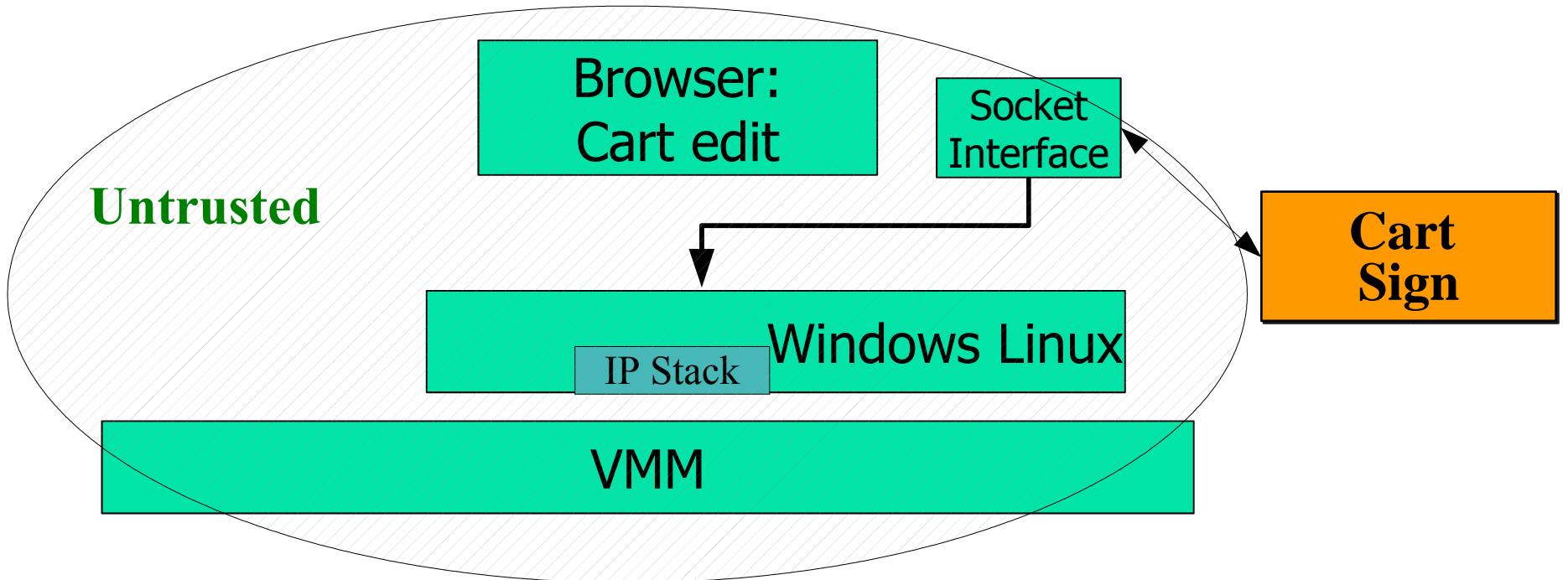


# Challenge: Untrusted VMM

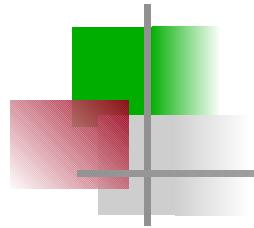


**Minimal Trusted Platform**

# Challenge: Untrusted VMM



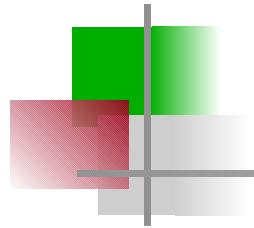
**Minimal Trusted Platform**



# Outline

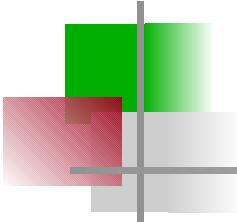
---

- propaganda
- some names
- security and system security objectives
- design principles
- architecture and components
  - use cases
  - components: current and future
- Nizza vs. Virtual Machine Monitors and other related work
- technical risks



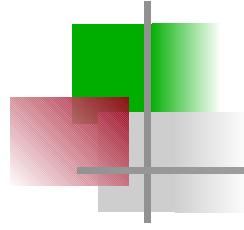
# Names

- TUD◆OS Technische Universität Dresden OS
  - DROPS Dresden Real-Time OS
  - Nizza Security Architecture
    - Micro-Sina A Nizza Application (use case)
  - L<sup>4</sup>Env set of servers and libraries
    - DOpE Window manager for DROPS and Nizza
  - L<sup>4</sup>Linux Linux kernel as user-level server
- L4 a micro-kernel interface
- L4/Fiasco, L4/Pistacchio: L4 implementations



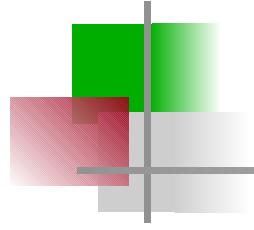
# Objectives: Security

- confidentiality  
no unauthorized access to information
- integrity  
no unauthorized, **unnoticed** modification of information
- recoverability  
no permanent damage to information
- availability  
timeliness of service



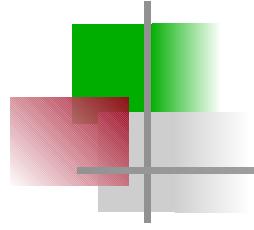
# Objectives: System Security

- Secure and unsecure applications



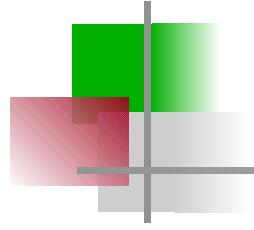
# Objectives: System Security

- Secure and unsecure applications
- Compatibility:
  - Legacy applications
  - Legacy/Fashionable Hardware



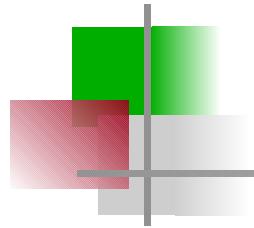
# Objectives: System Security

- Secure and unsecure applications
- Compatibility
- Flexible sandboxing



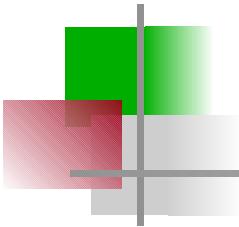
# Objectives: System Security

- Secure and unsecure applications
- Compatibility
- Flexible sandboxing
- Resource Control



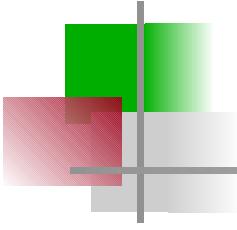
# Objectives: System Security

- Secure and unsecure applications
- Compatibility
- Flexible sandboxing
- Resource Control
- Small TCB:  
complexity acceptable, if for a small group
  - Each member fully understands interaction of all components.
  - Each component is understood by one member.



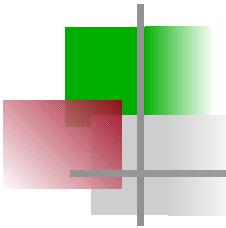
# Principles

- Trusted Computing Base: per application
- platform:
  - small set of small components (servers, ...)
  - small interfaces
  - select components of platform per application



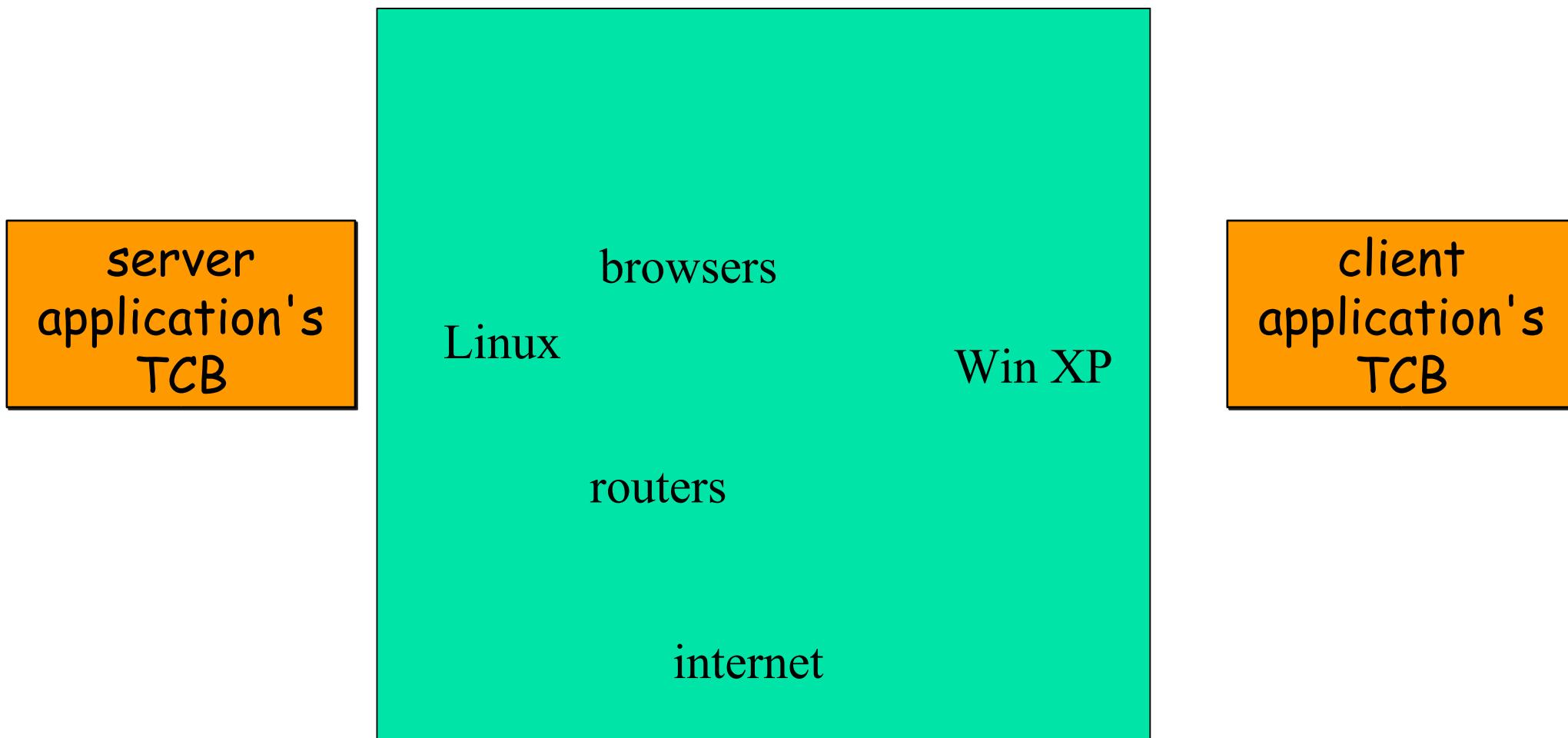
# Principles, continued

- split applications and services:  
sensitive part in/on trusted platform  
and *other* part
- reuse legacy for *other* part  
-> trusted wrappers / tunneling



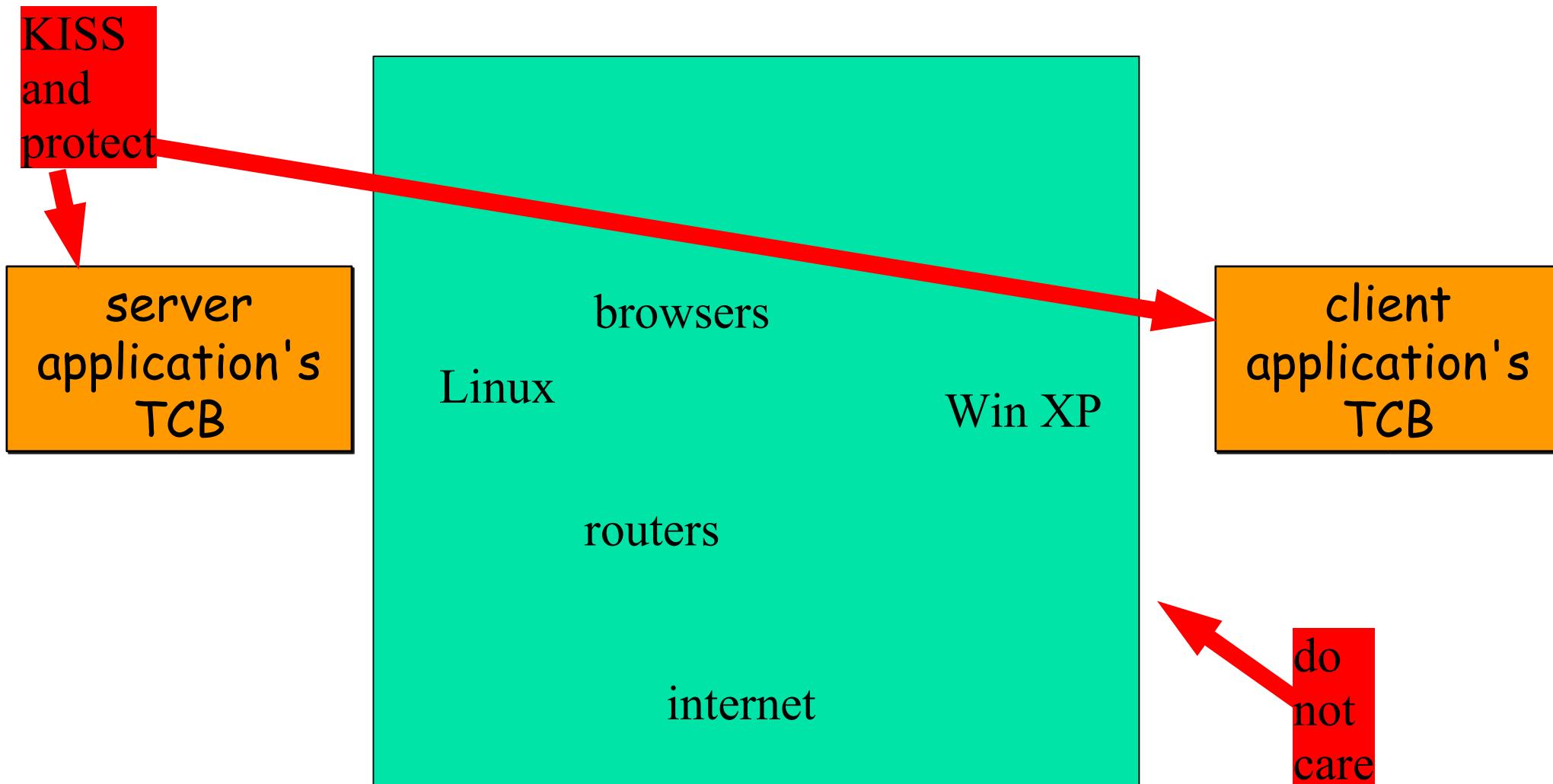
# Principles

- push end-to-end argument to the extreme



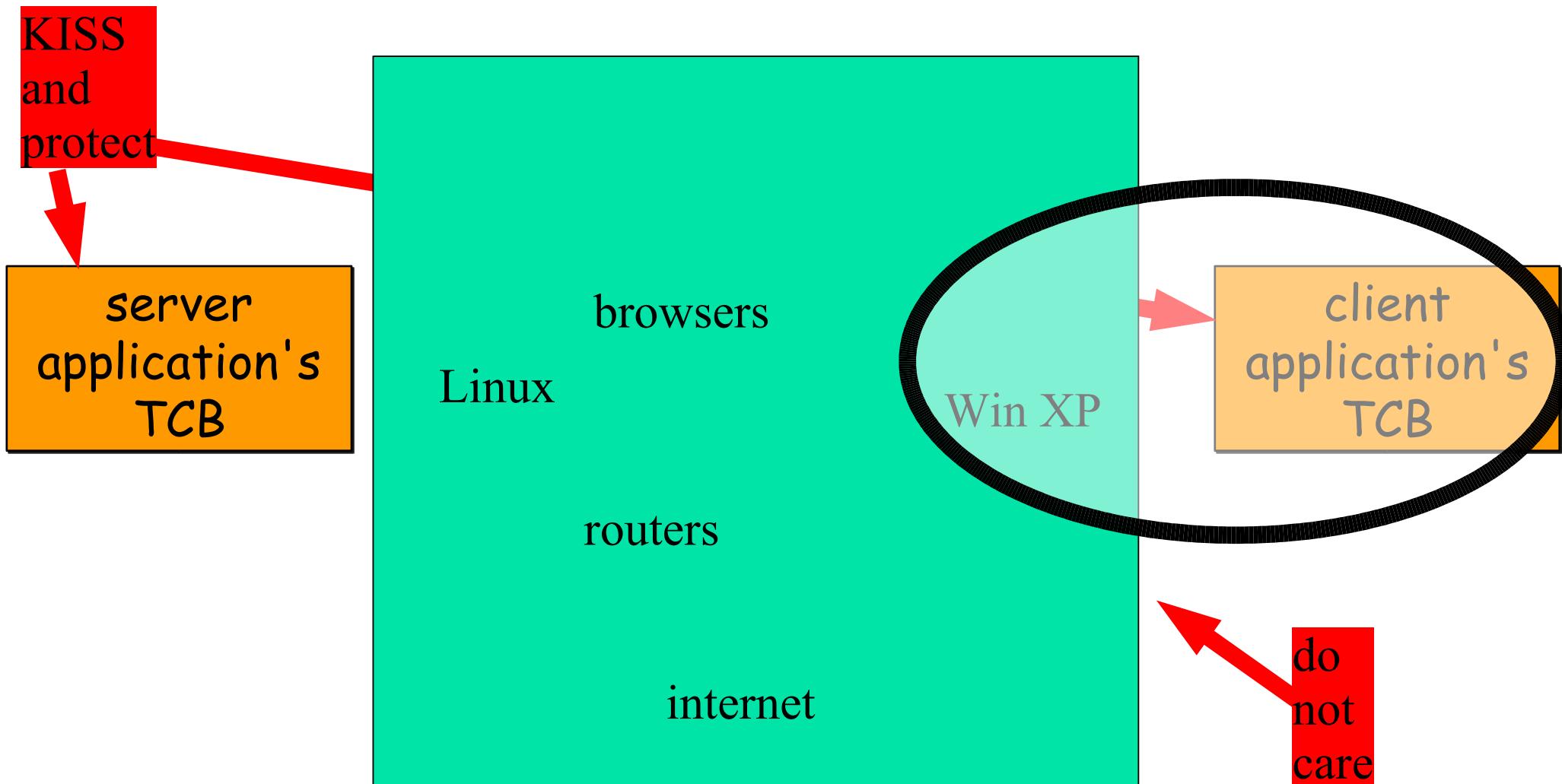
# Principles

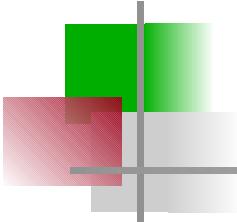
- push end-to-end argument to the extreme



# Principles

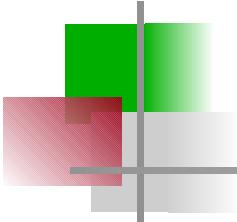
- push end-to-end argument to the extreme





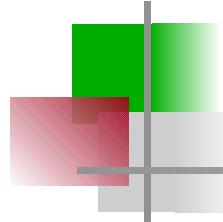
# Principles and Techniques

- micro kernel:
  - separates legacy from sensitive partitions
  - separates components of small platform
  - provides mechanisms for access control  
mediates communication
- contract-based access control
- secure booting / attestation +  
trusted path to user

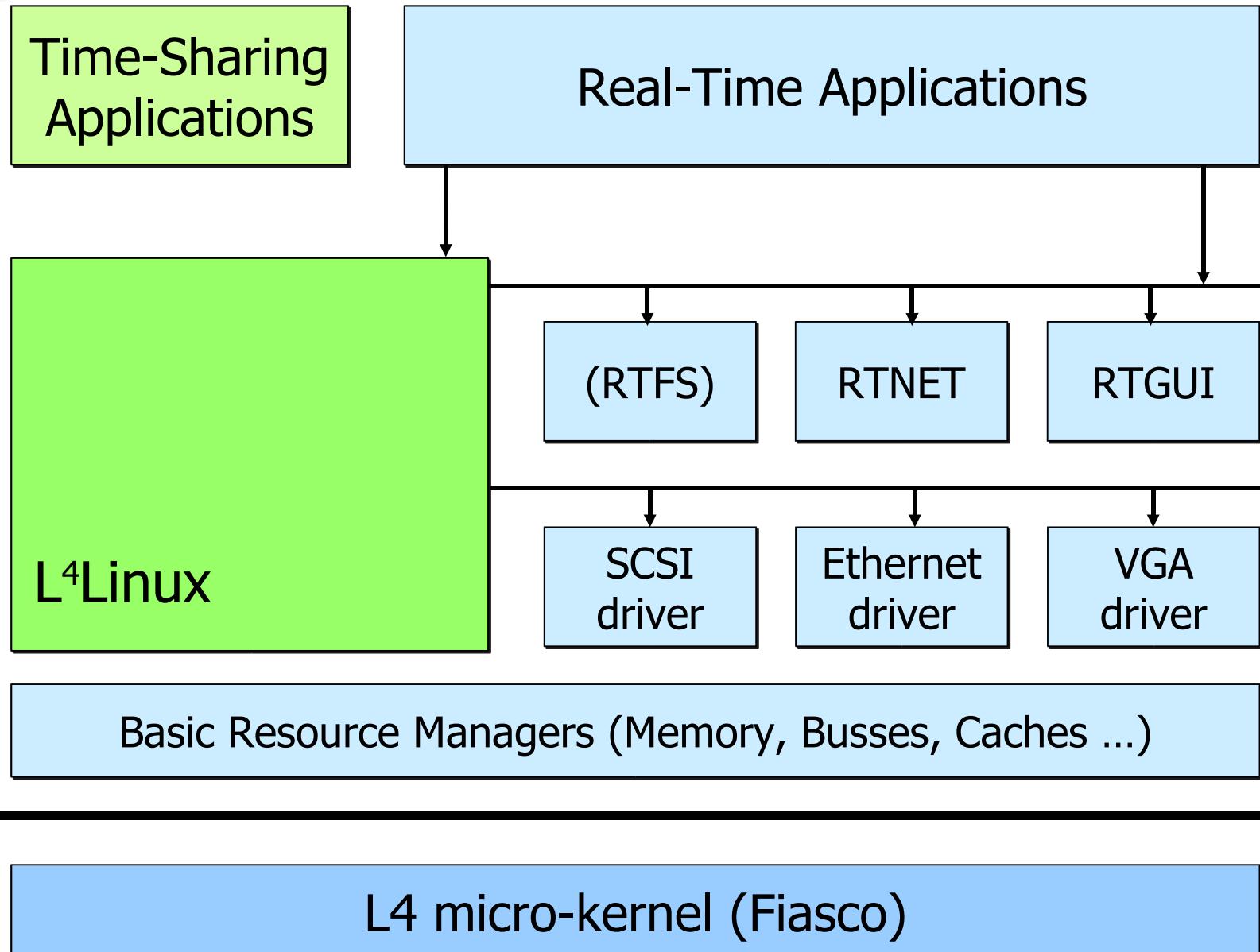


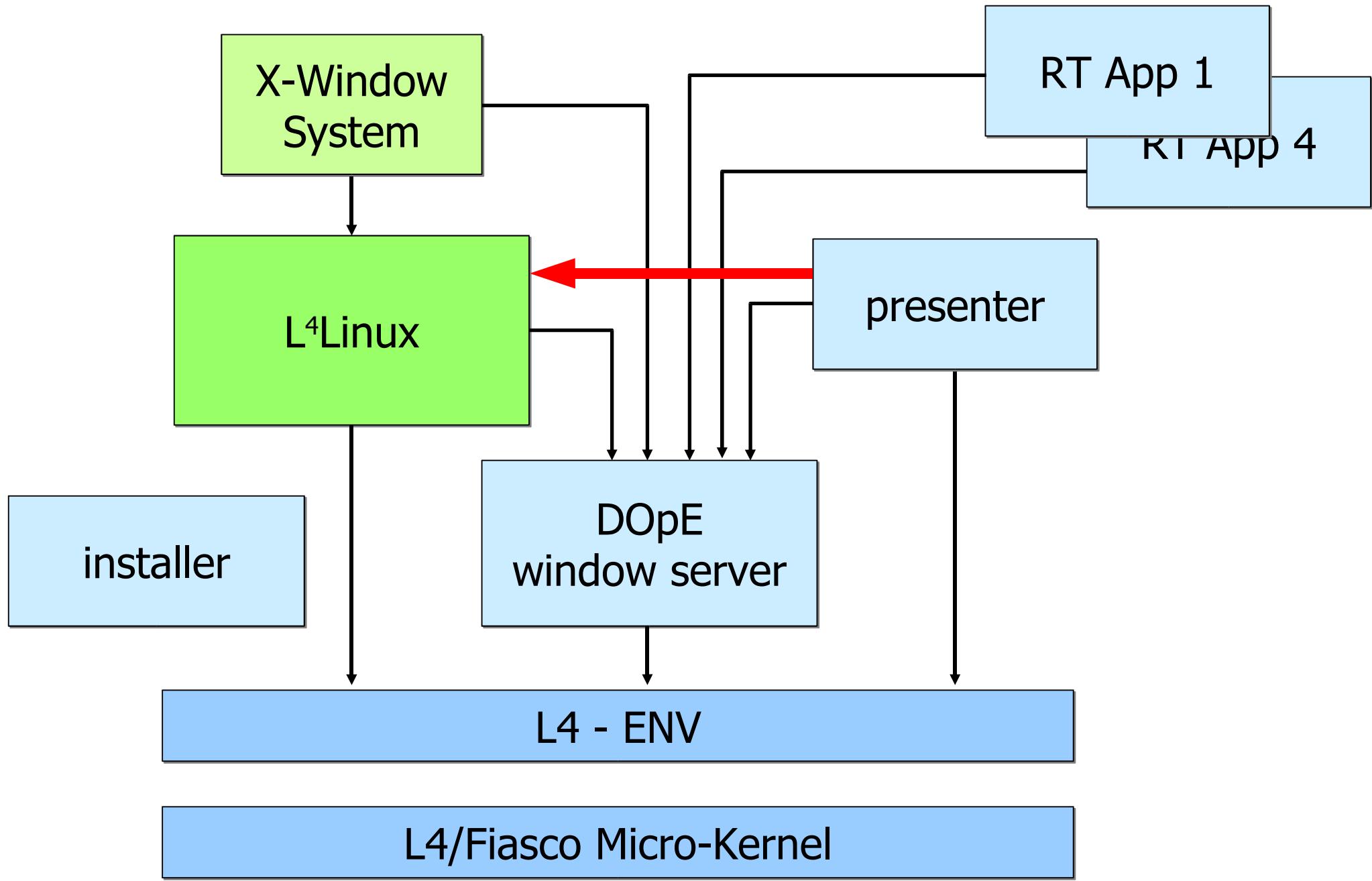
# Architecture + Components

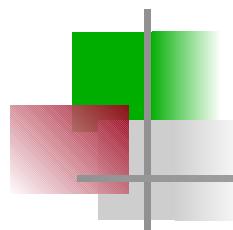
- the origin: Dresden Real-time OPerating System
- Nizza architecture
- use cases
- some Nizza components



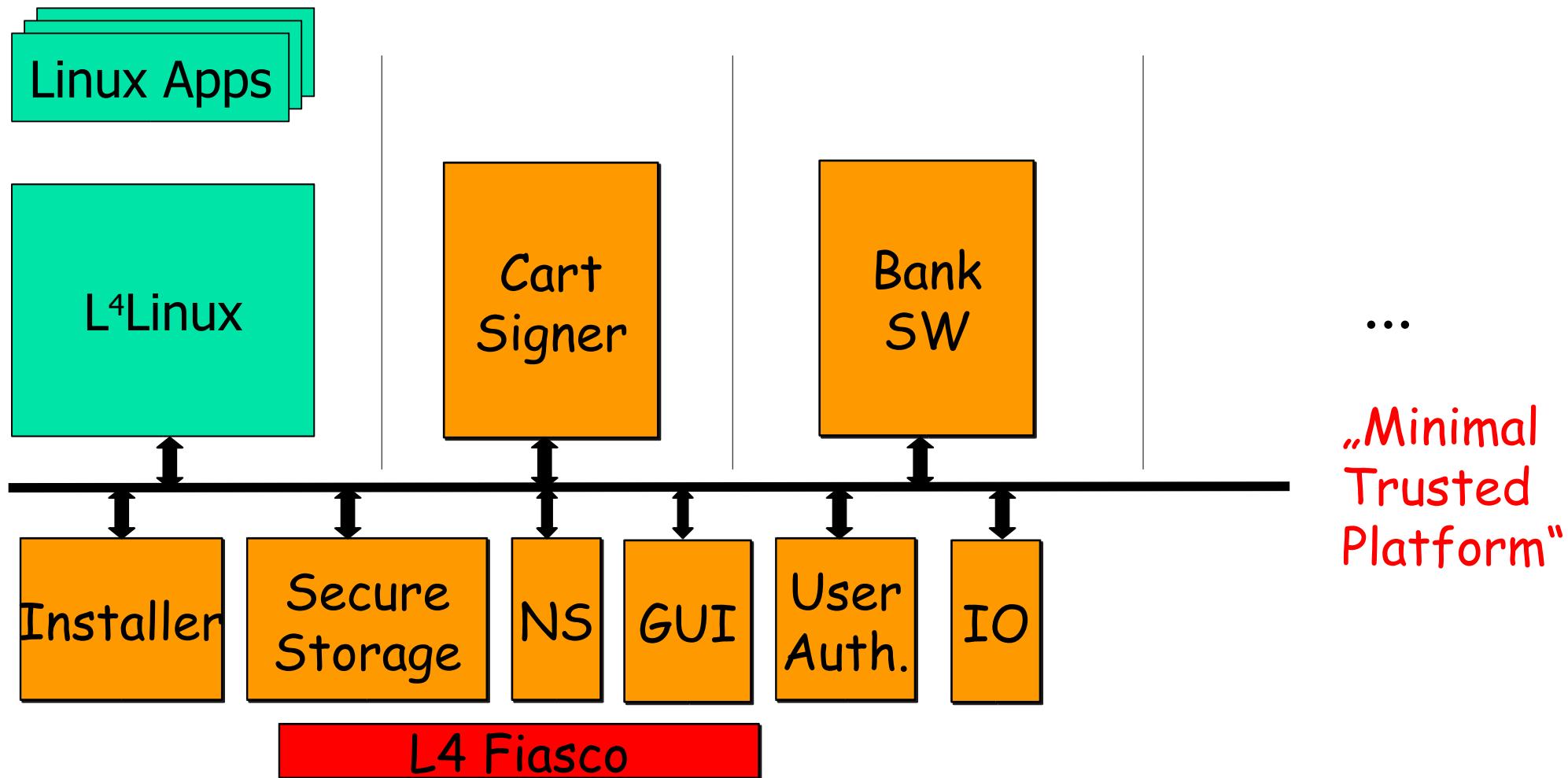
# Starting Point: DResden Real-Time OS



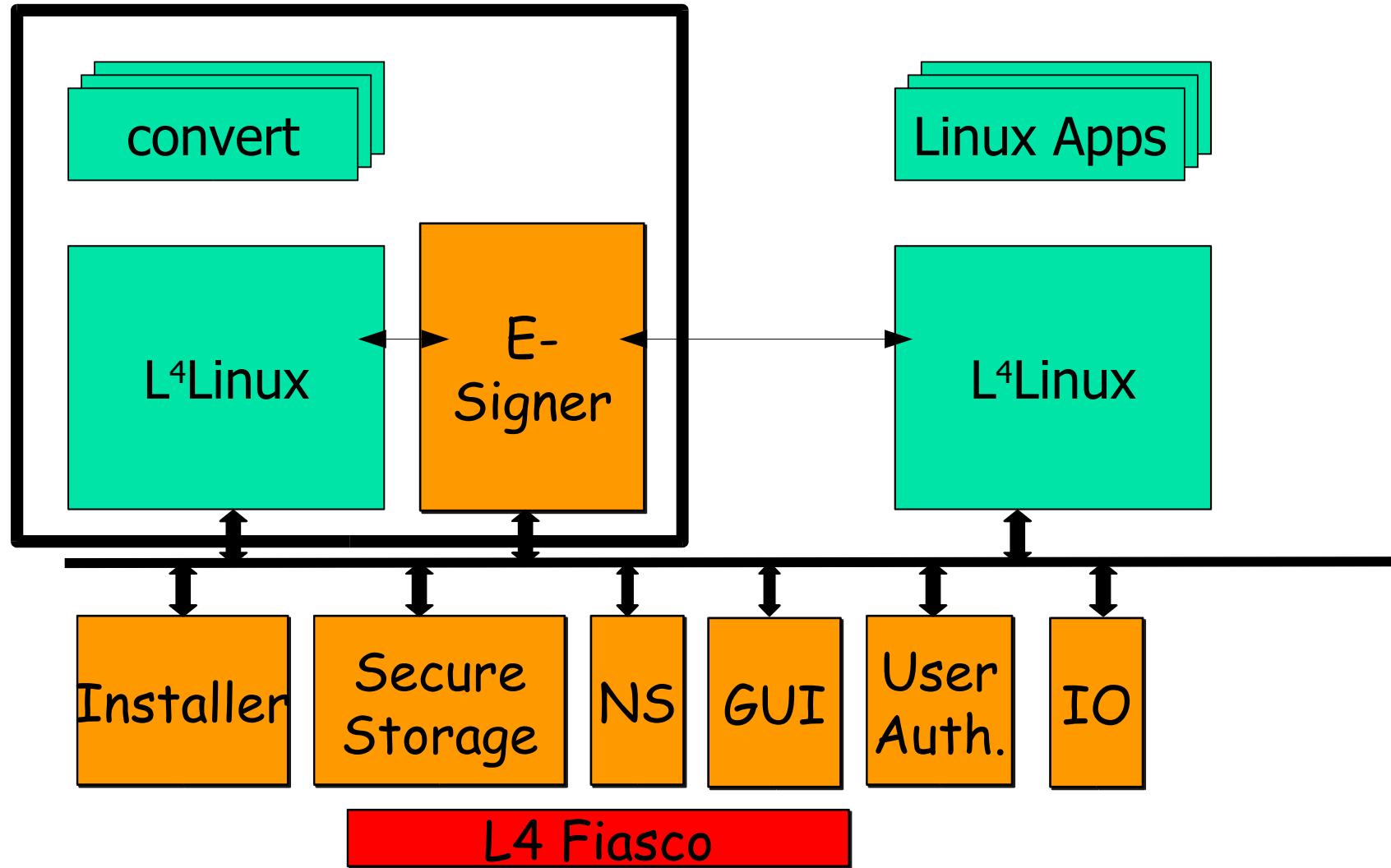


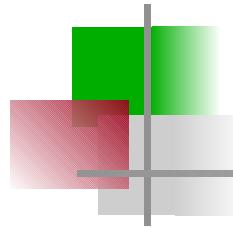


# NIZZA Architecture

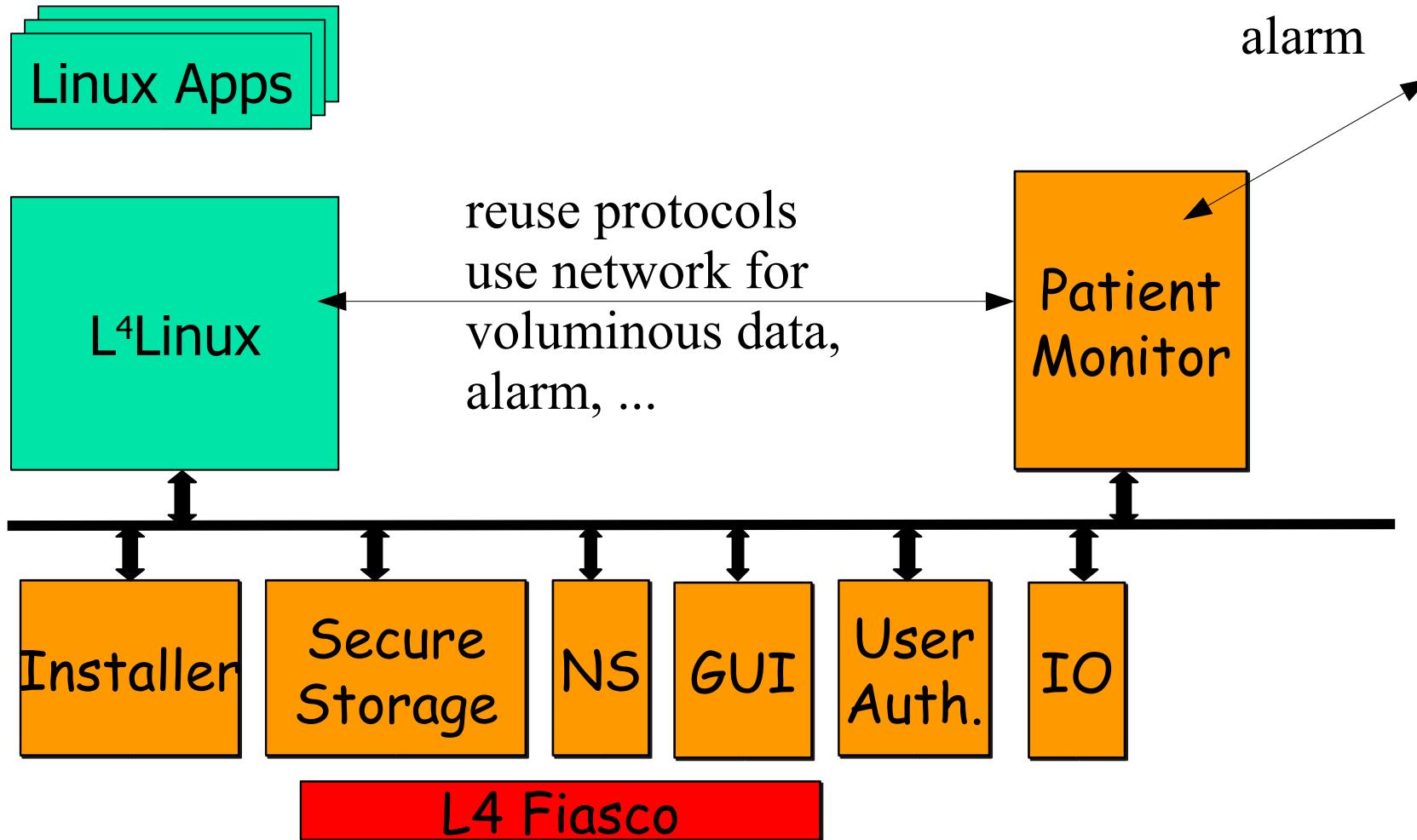


# More Use Cases

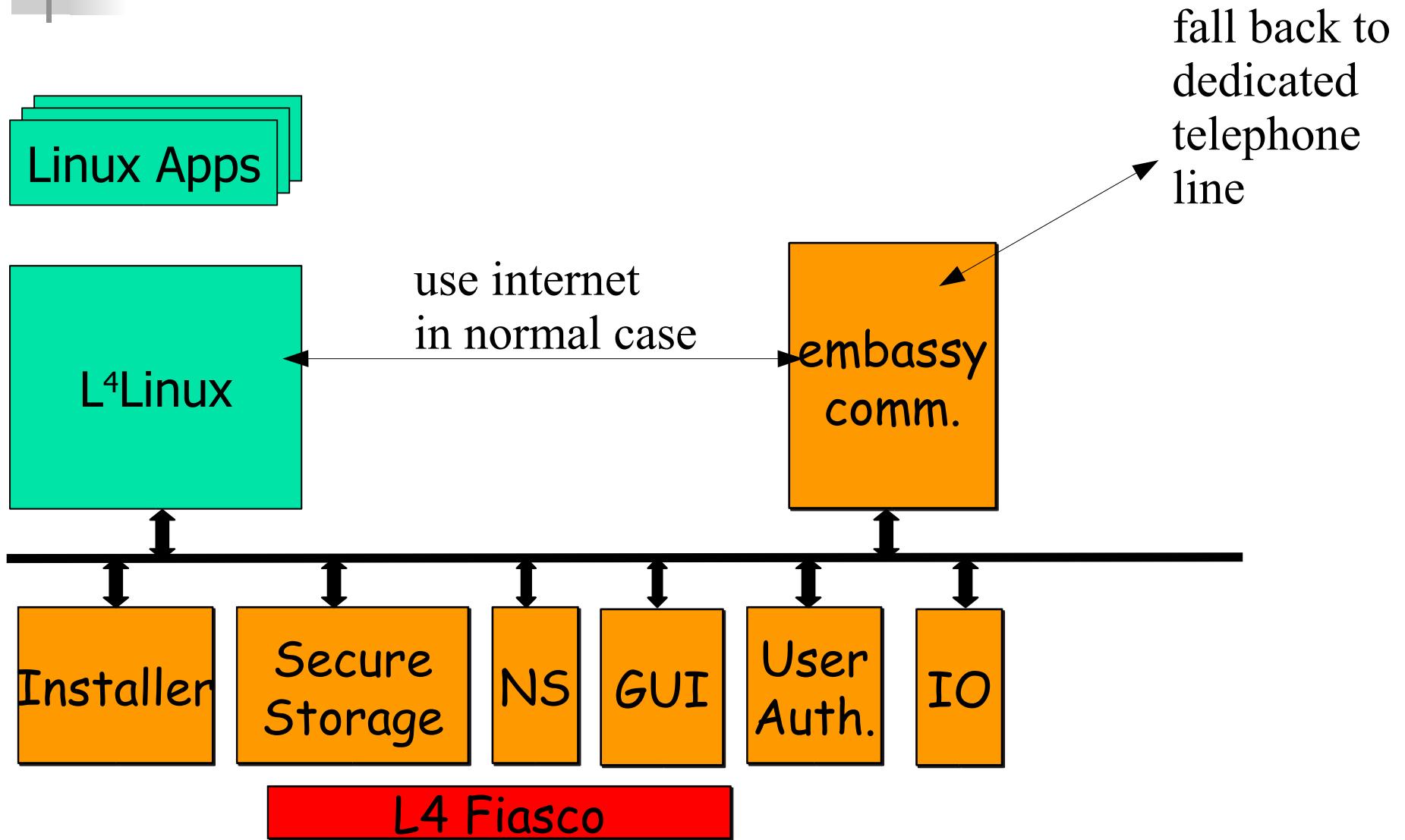




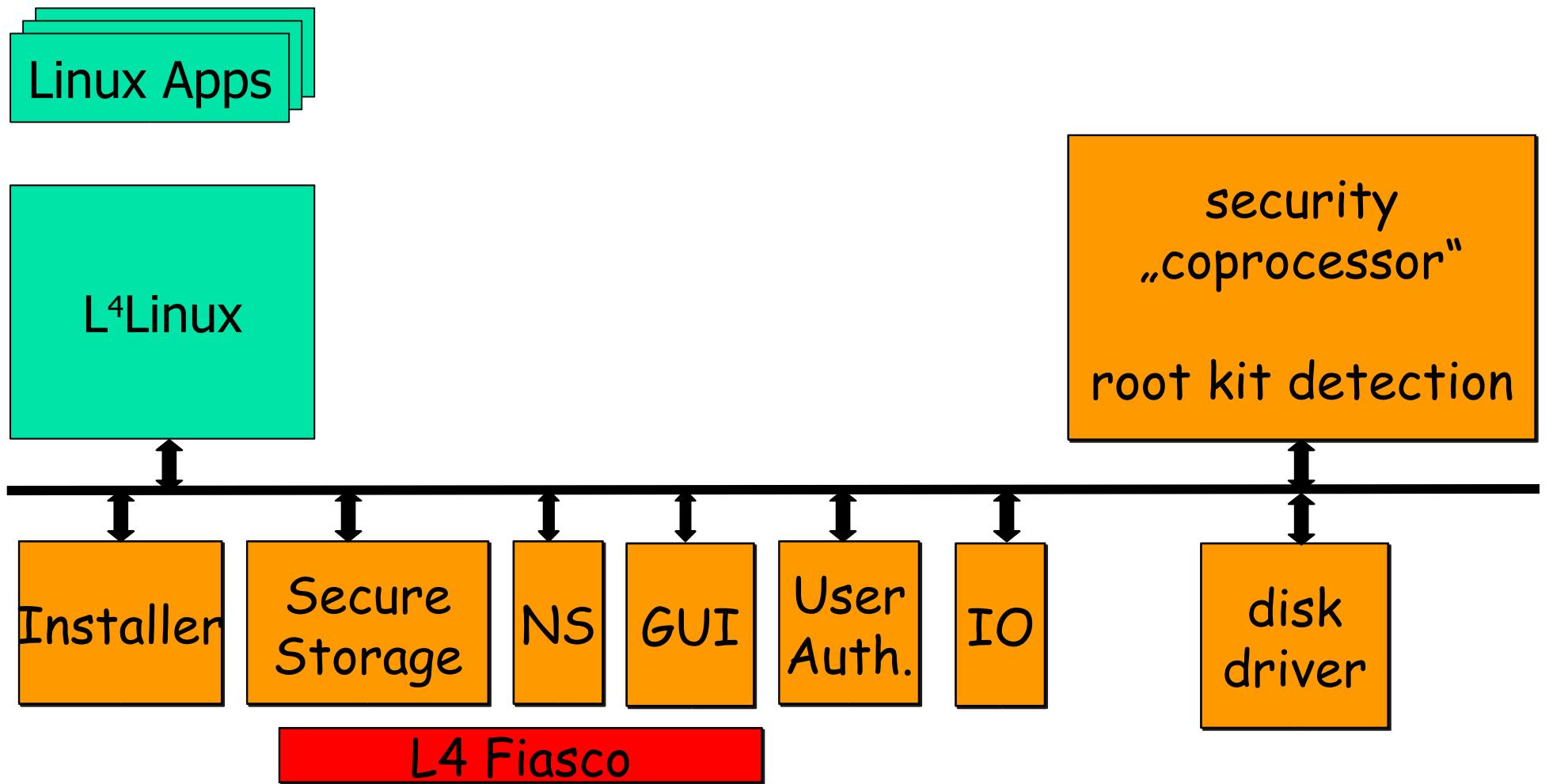
# More Use Cases



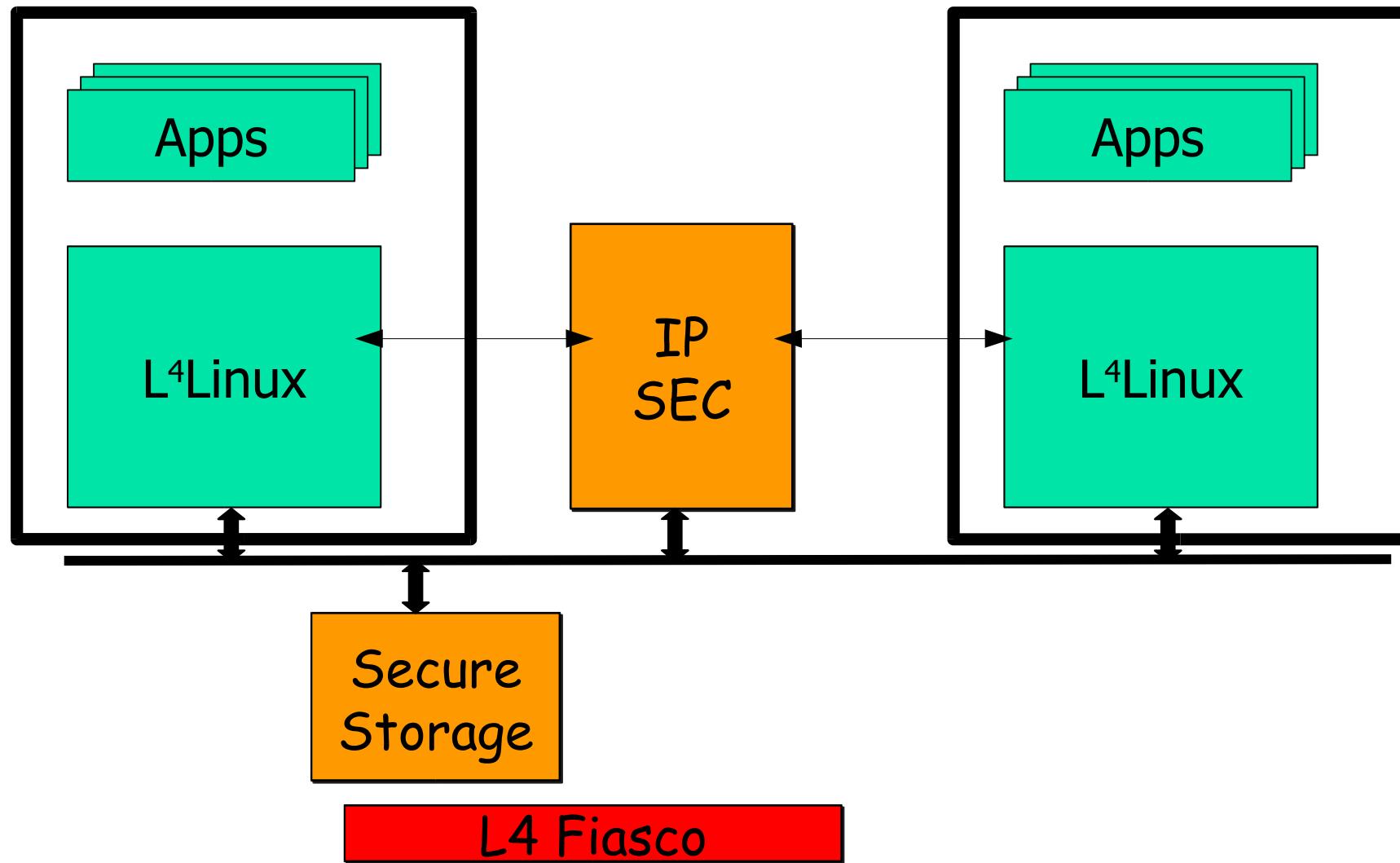
# More Use Cases

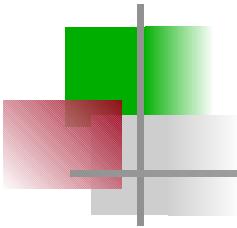


# More Use Cases



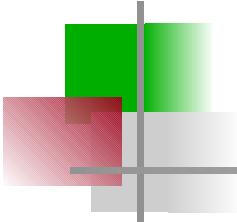
# More Use Cases





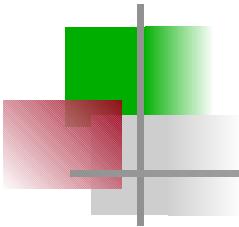
# Architecture + Components

- the origin: Dresden Real-time OPerating System
- Nizza architecture
- use cases
- **some Nizza components**
  - L4 micro kernel: present and future
  - L<sup>4</sup>Linux: encapsulation and reuse
  - secure booting + trusted path
  - secure storage with small TCB (future)



# L4 Micro-Kernel: evolution

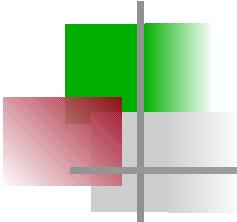
- the original:  
address spaces, threads, IPC
- L4/Fiasco-RT: real-time:  
periodic threads, fine grained scheduling support, etc
- **L4/Fiasco-X.e:**  
**unified access control + kernel resource management**
- L4-Next: virtualization support



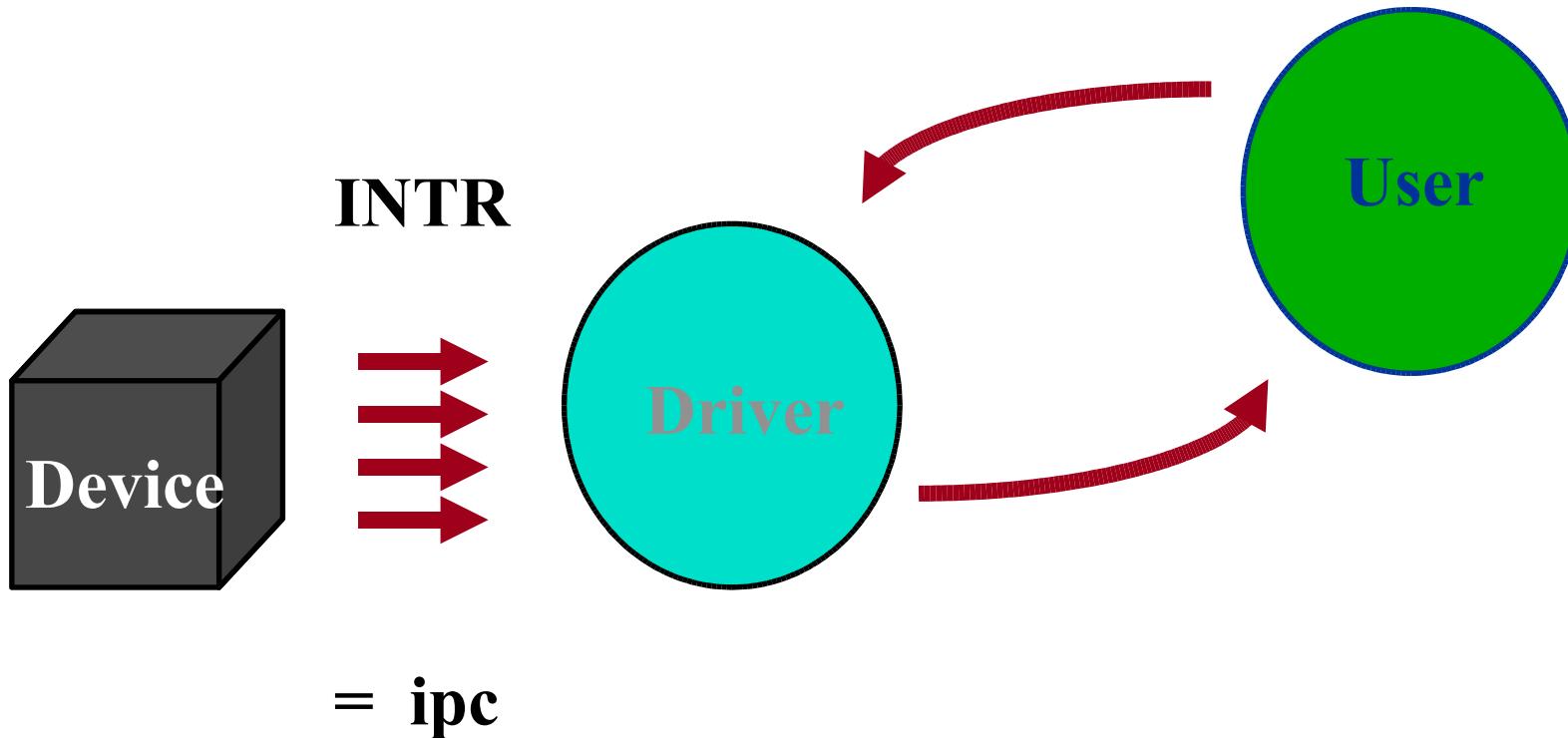
# L4 micro kernel (Jochen Liedtke)

## fundamental abstractions

- address spaces (separation)
- threads
- inter process communication (IPC)
  - explicit
  - interrupts
  - faults and mappings

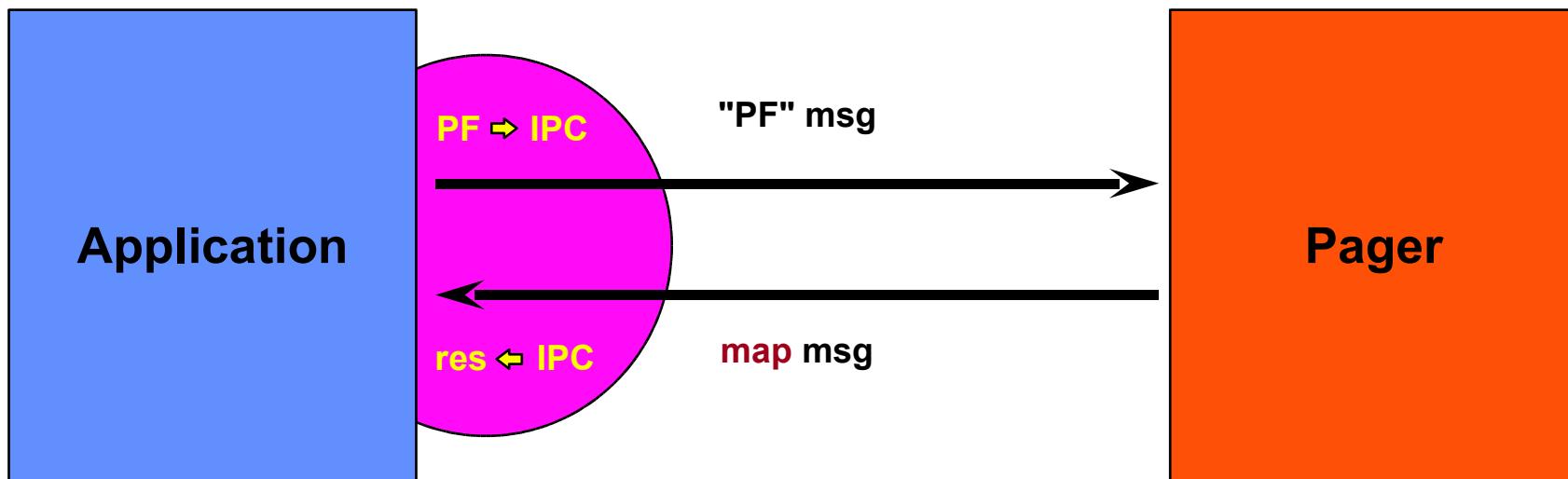


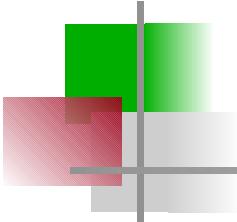
# Drivers at User Level



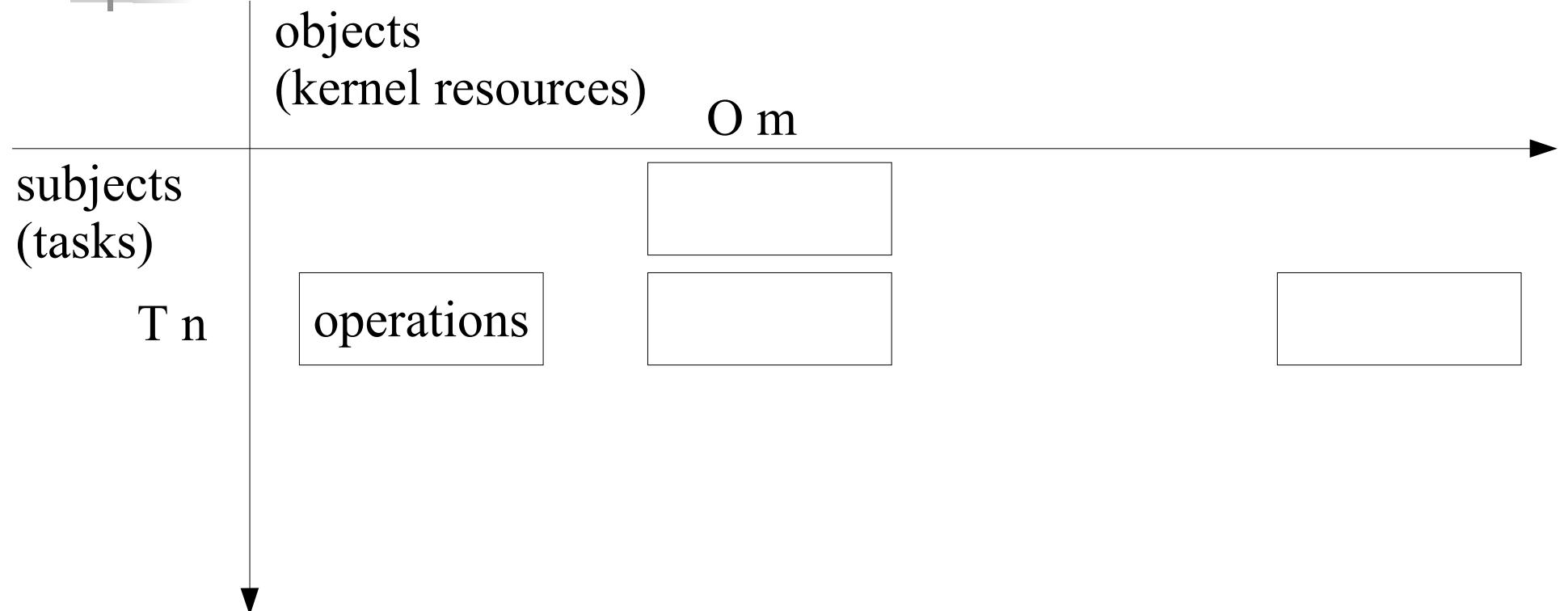
- **IO ports:** part of the user address space
- **interrupts:** messages from hardware

# Memory Pages

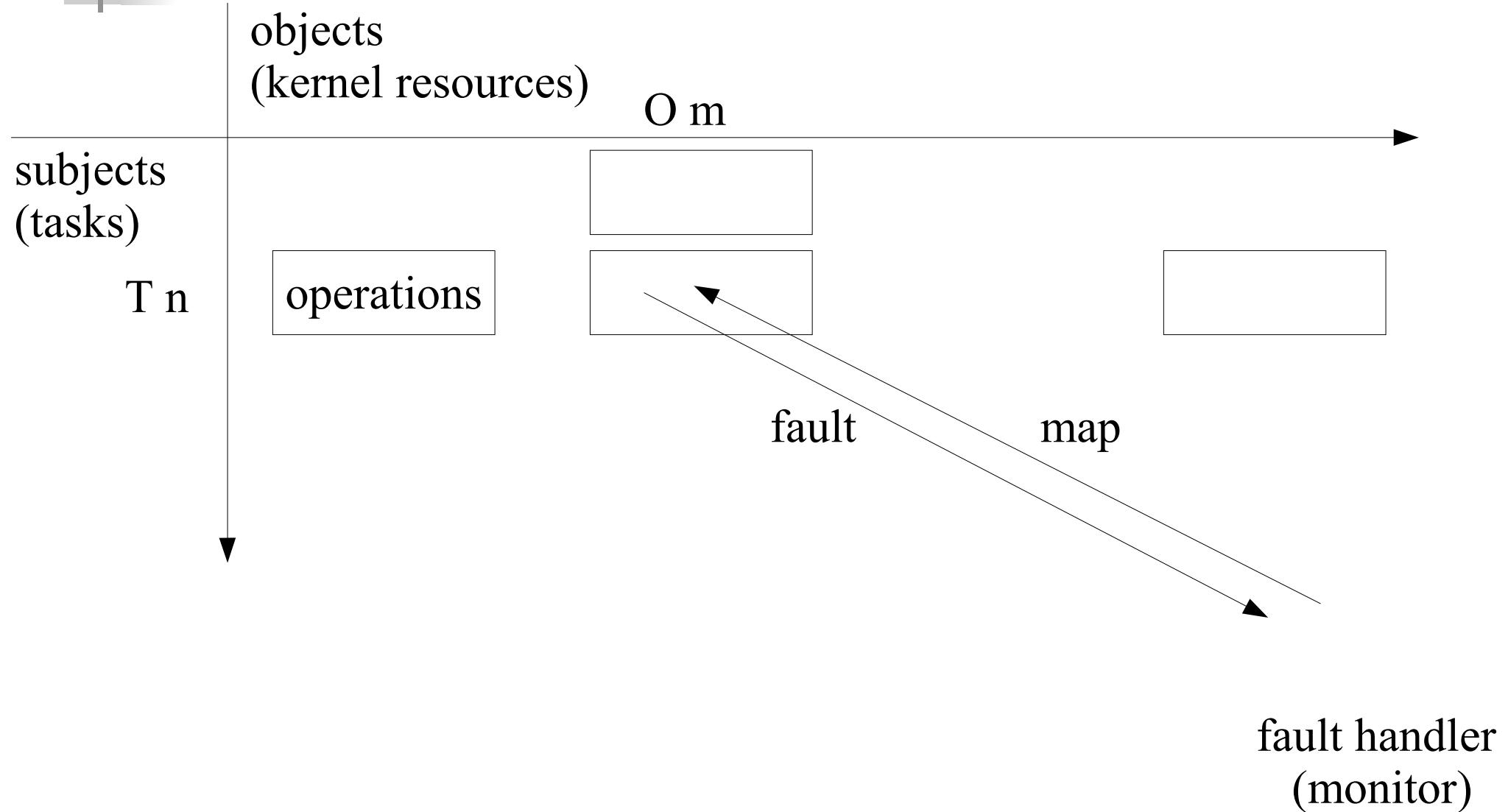




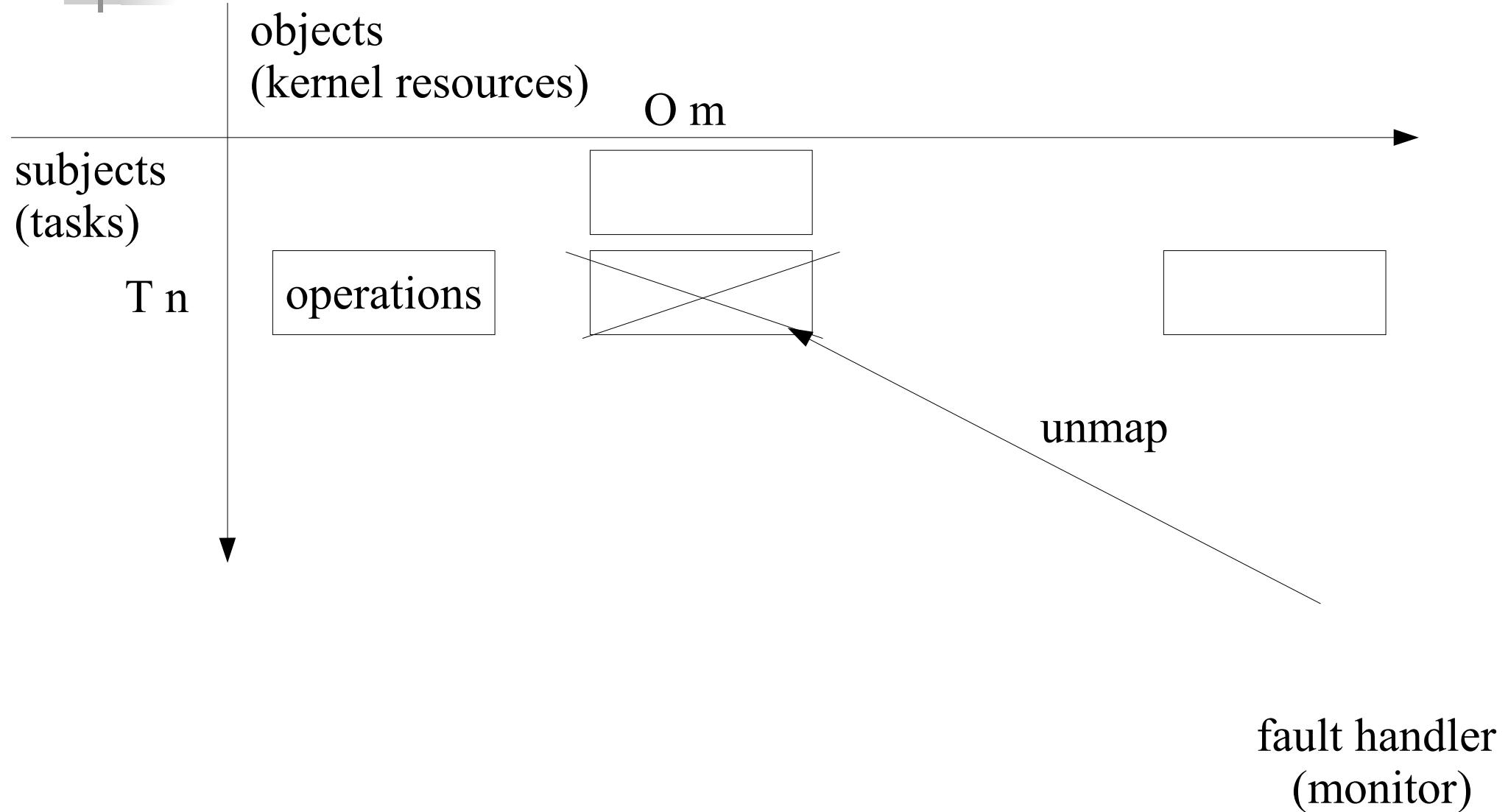
# L4 X.e: Unifying Access Control

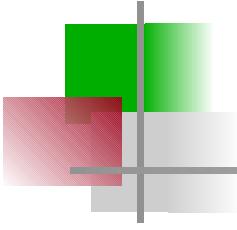


# L4 X.e: Unifying Access Control



# L4 X.e: Unifying Access Control

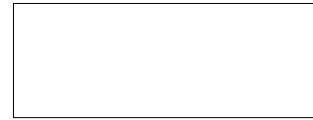


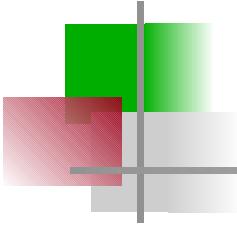


# L4 X.e: Unifying Access Control

task: T n

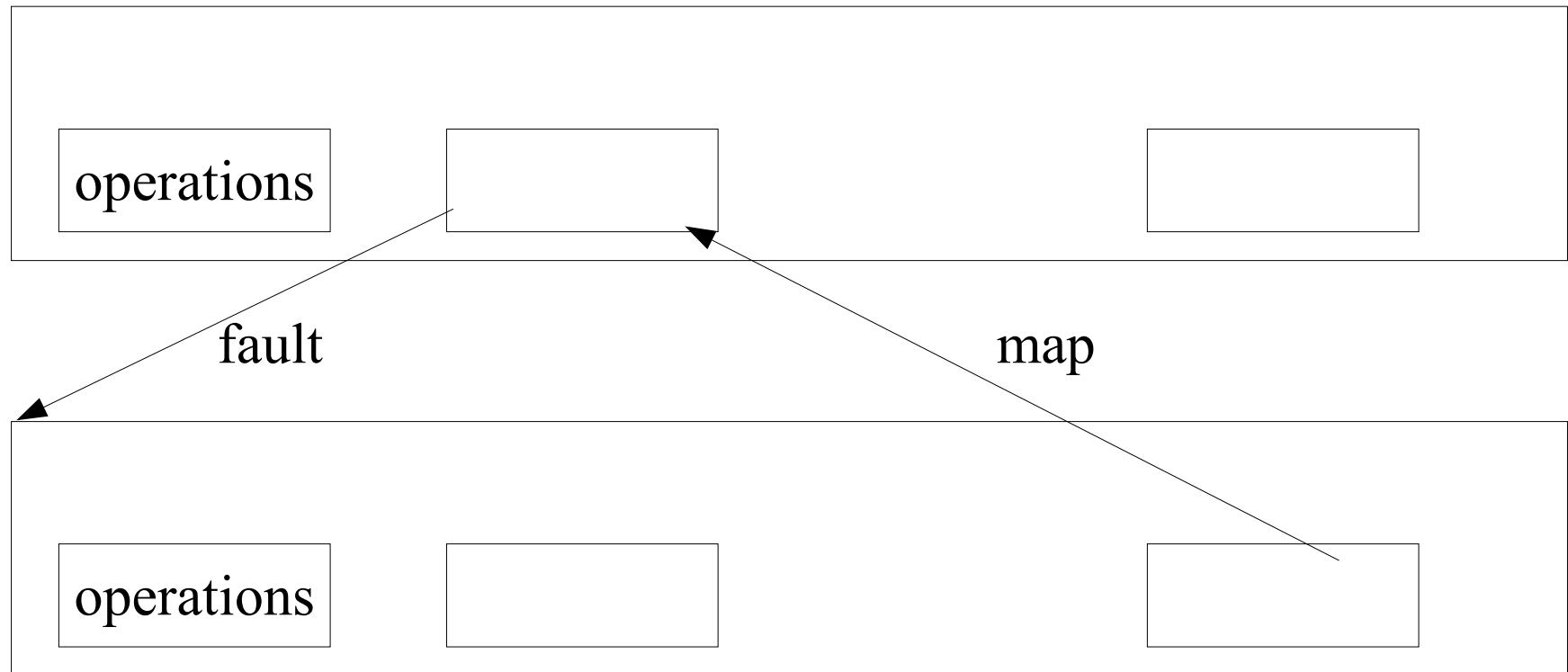
operations



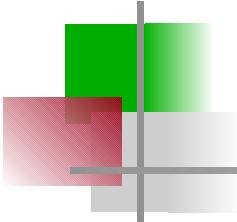


# L4 X.e: Unifying Access Control

task: T n

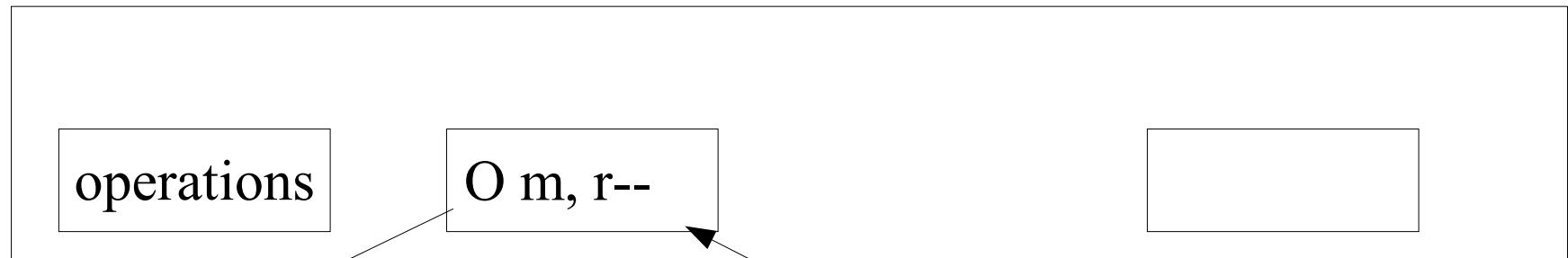


fault handler  
(monitor)



# L4 X.e: Unifying Access Control

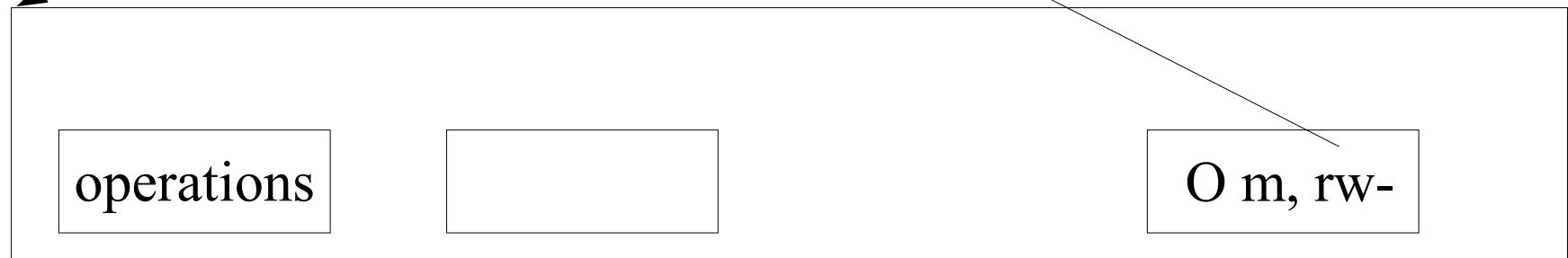
task: T n



fault

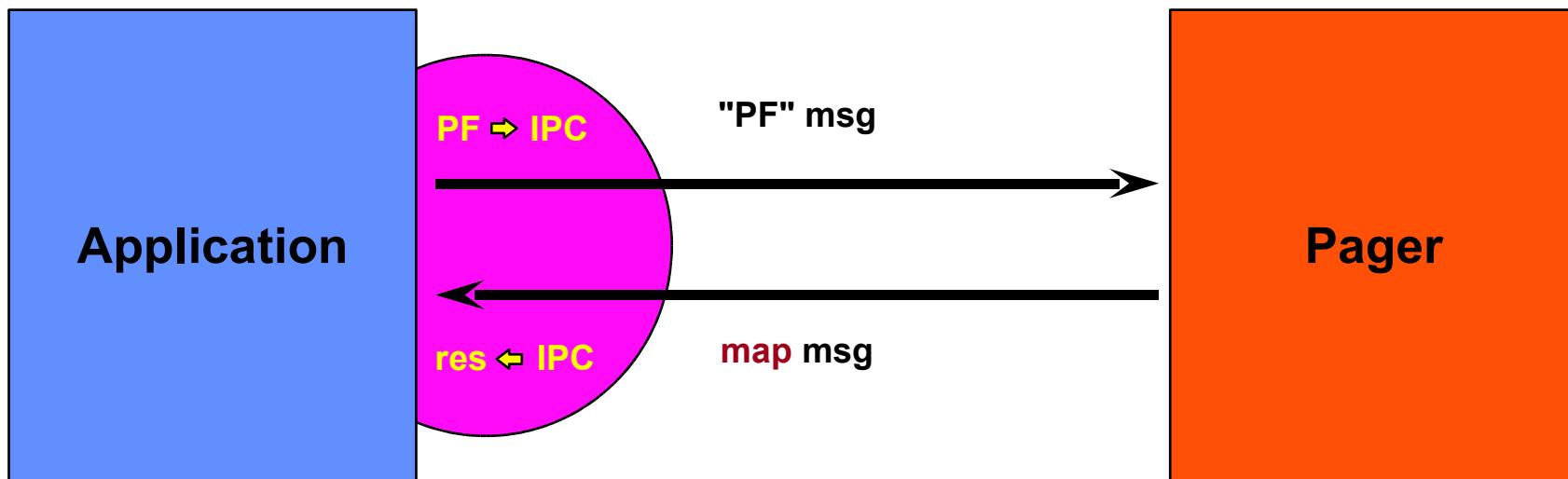
map r-x

task: T x

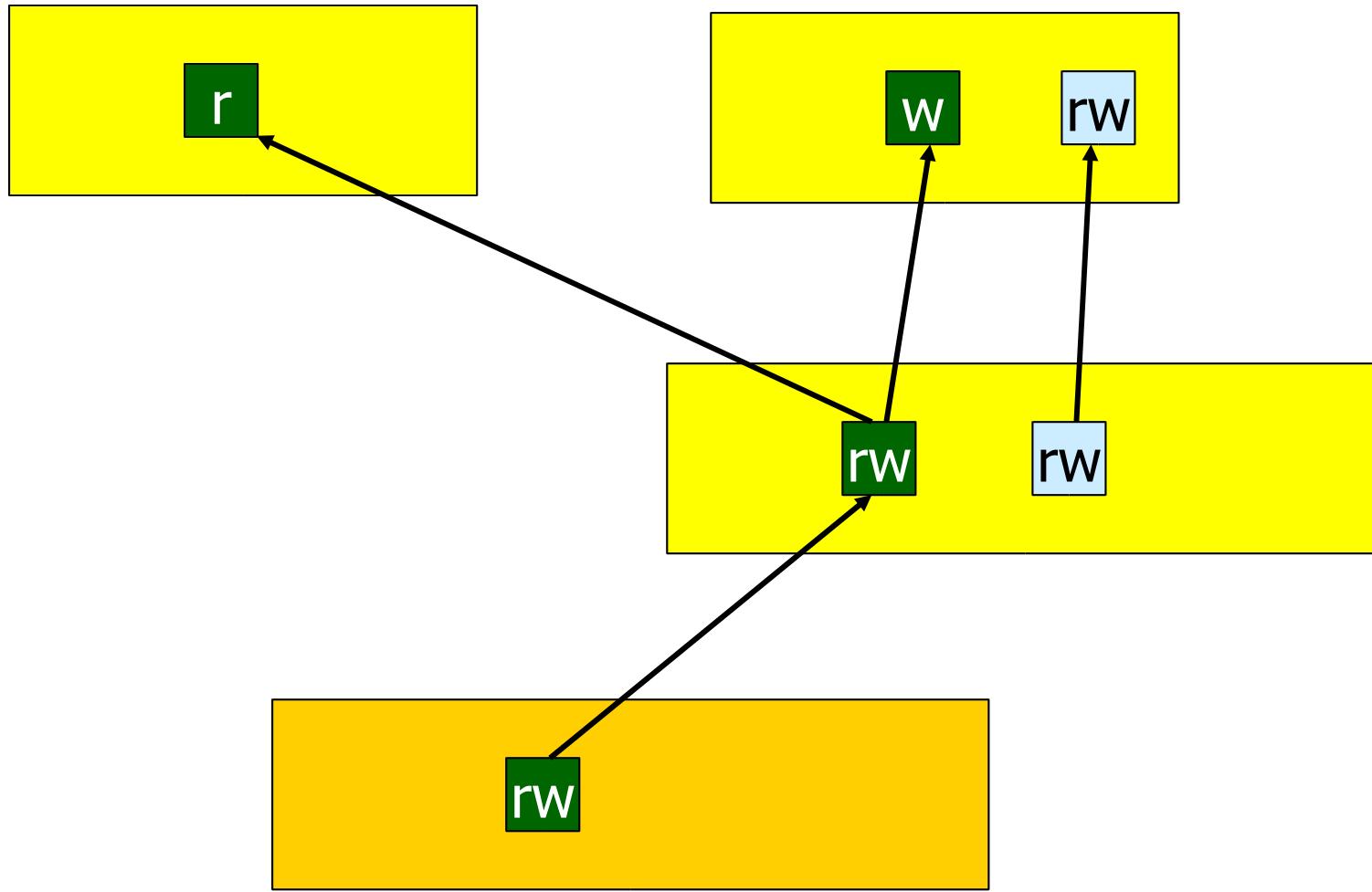


fault handler  
(monitor)

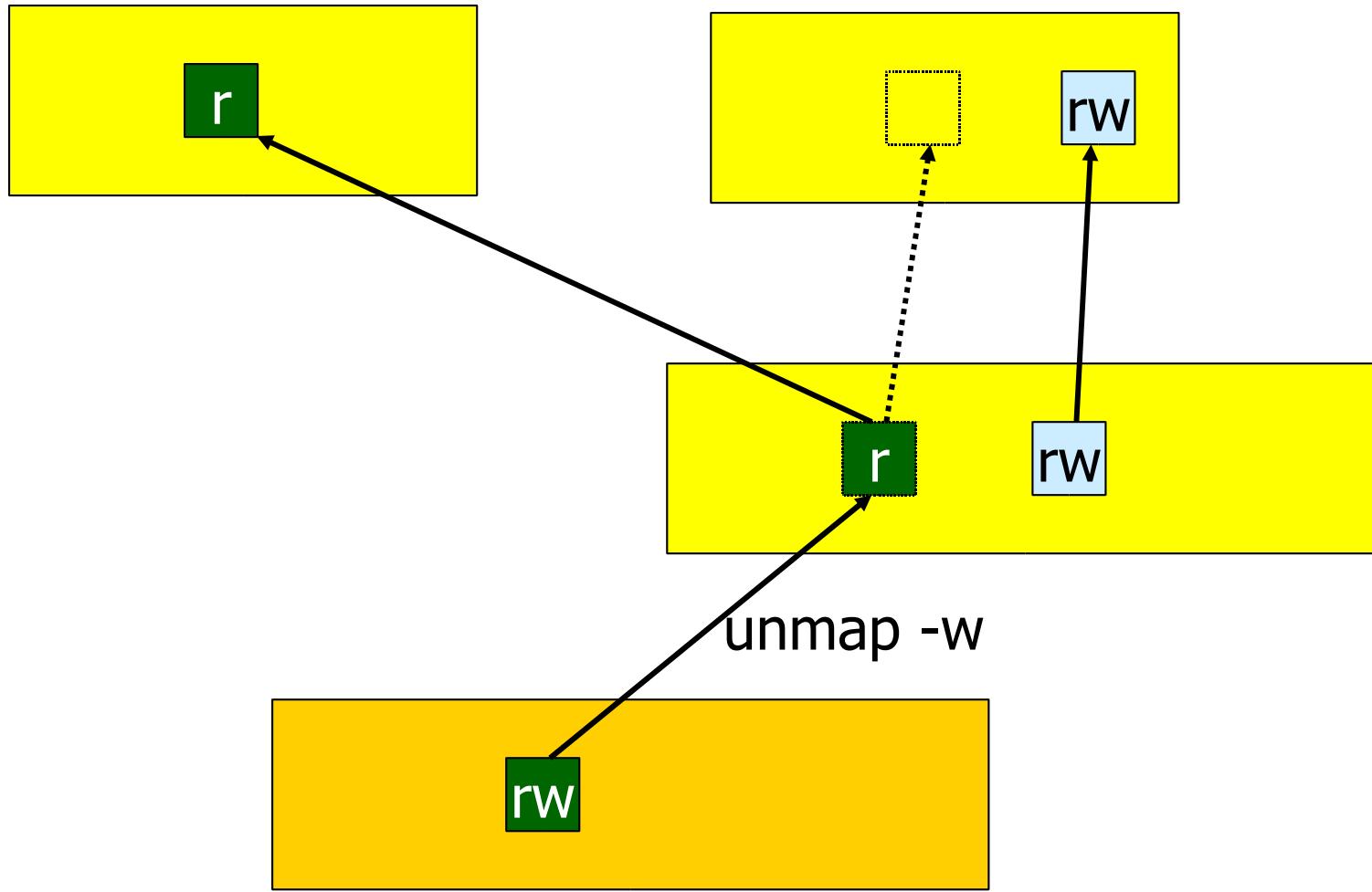
# Memory Pages



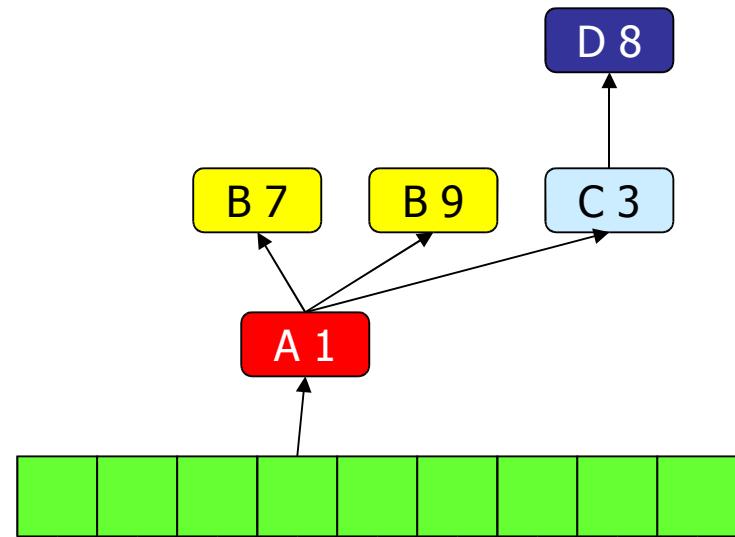
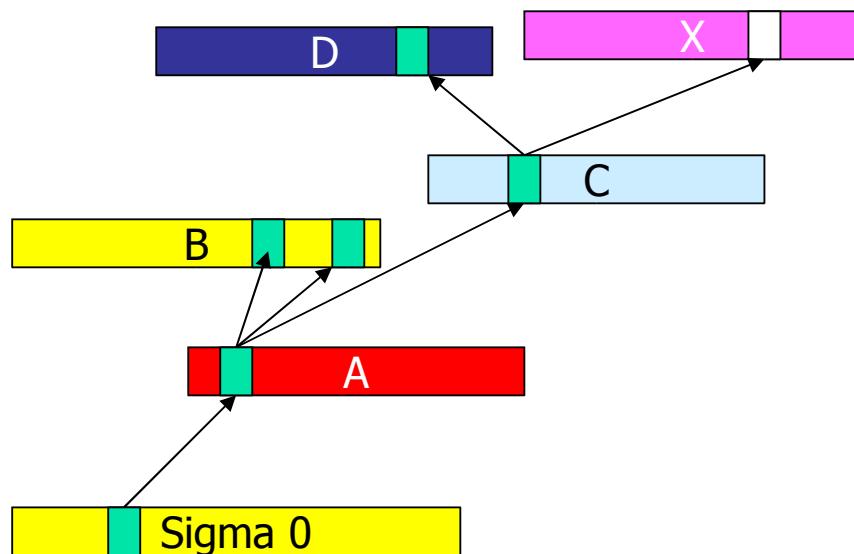
# Fault, Map, Unmap



# Fault, Map, Unmap

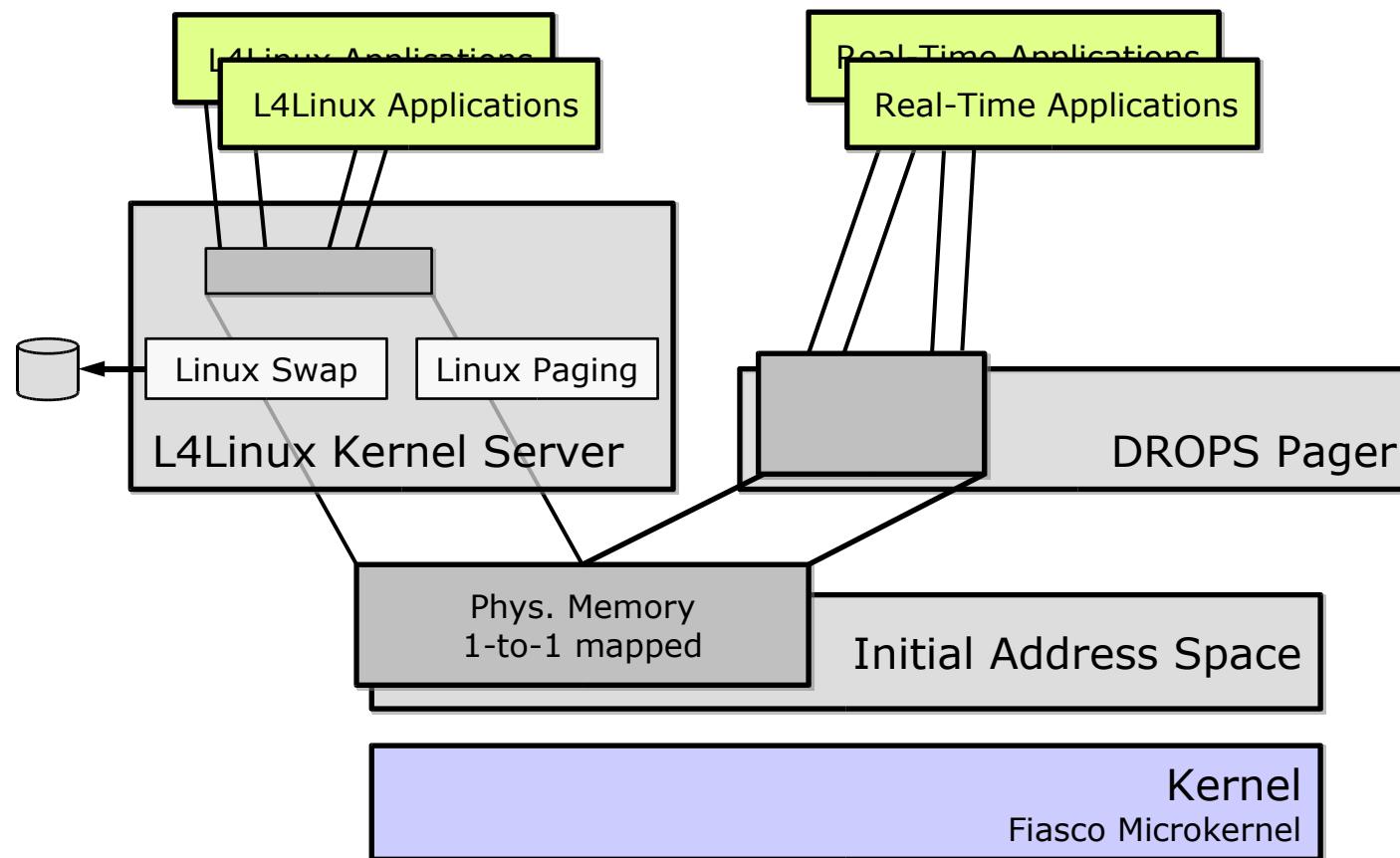


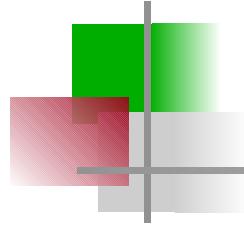
# Map Trees



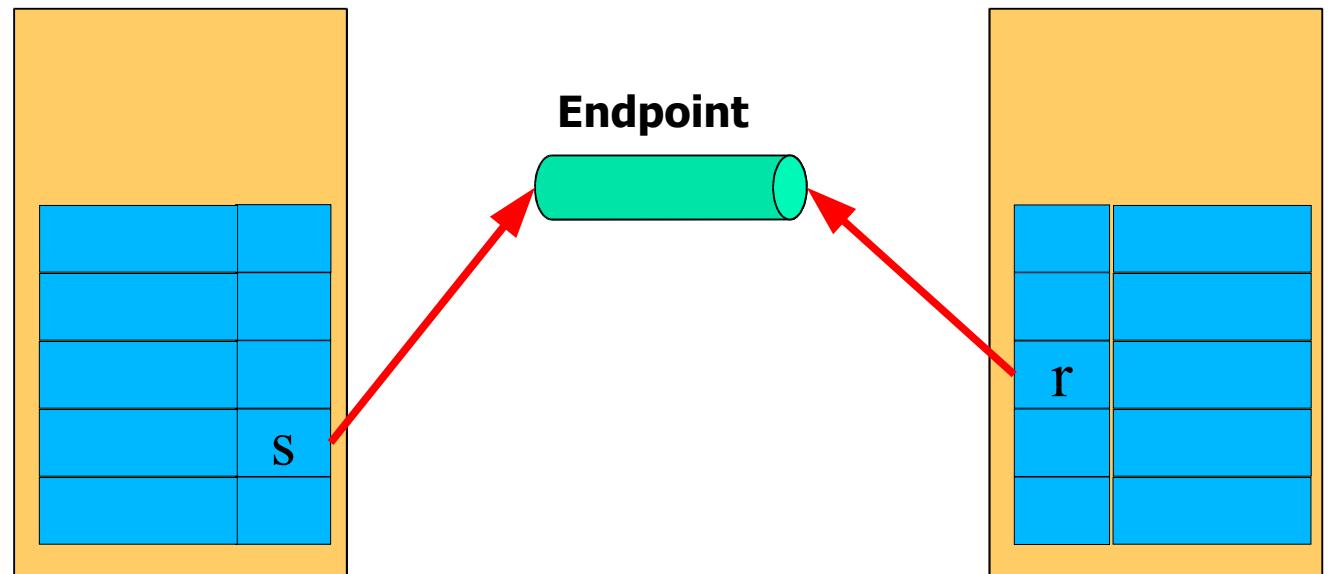
Tiles (physical memory frames)

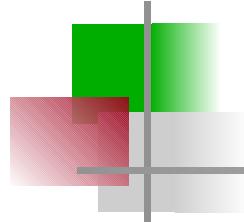
# Pager Example



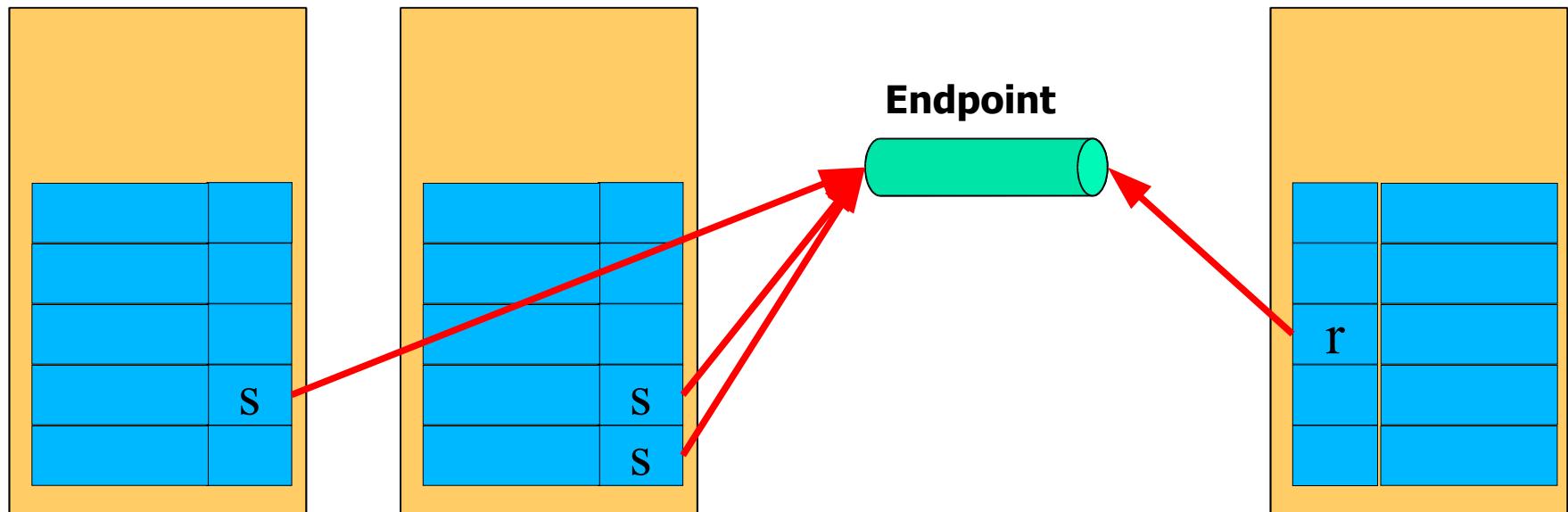


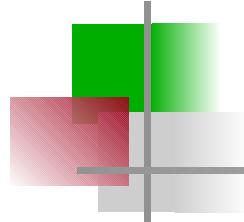
# Communication



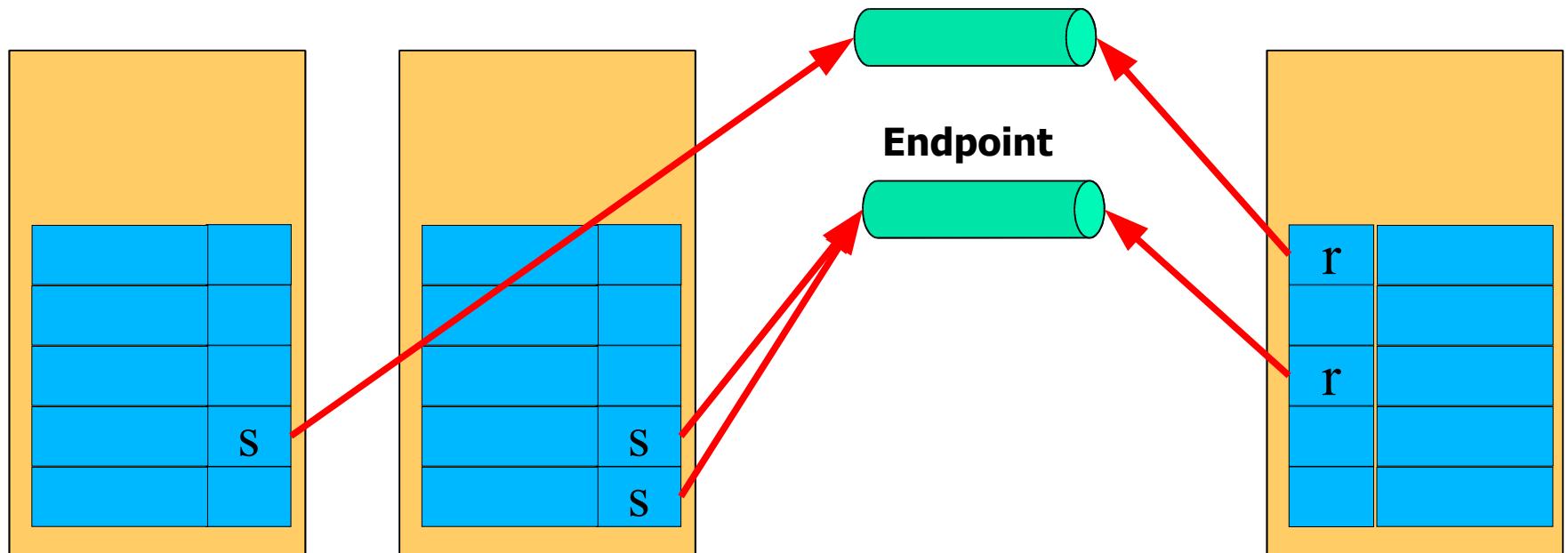


# Communication

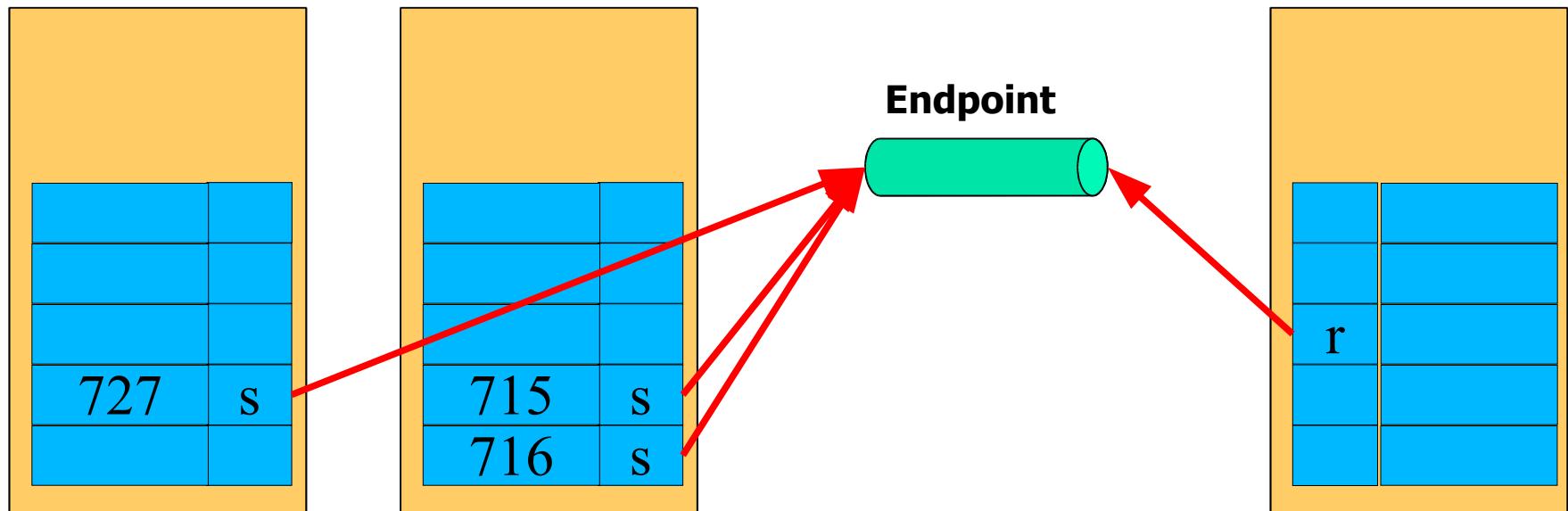


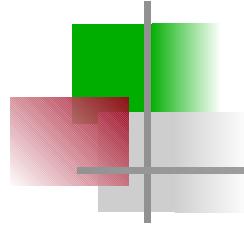


# Communication

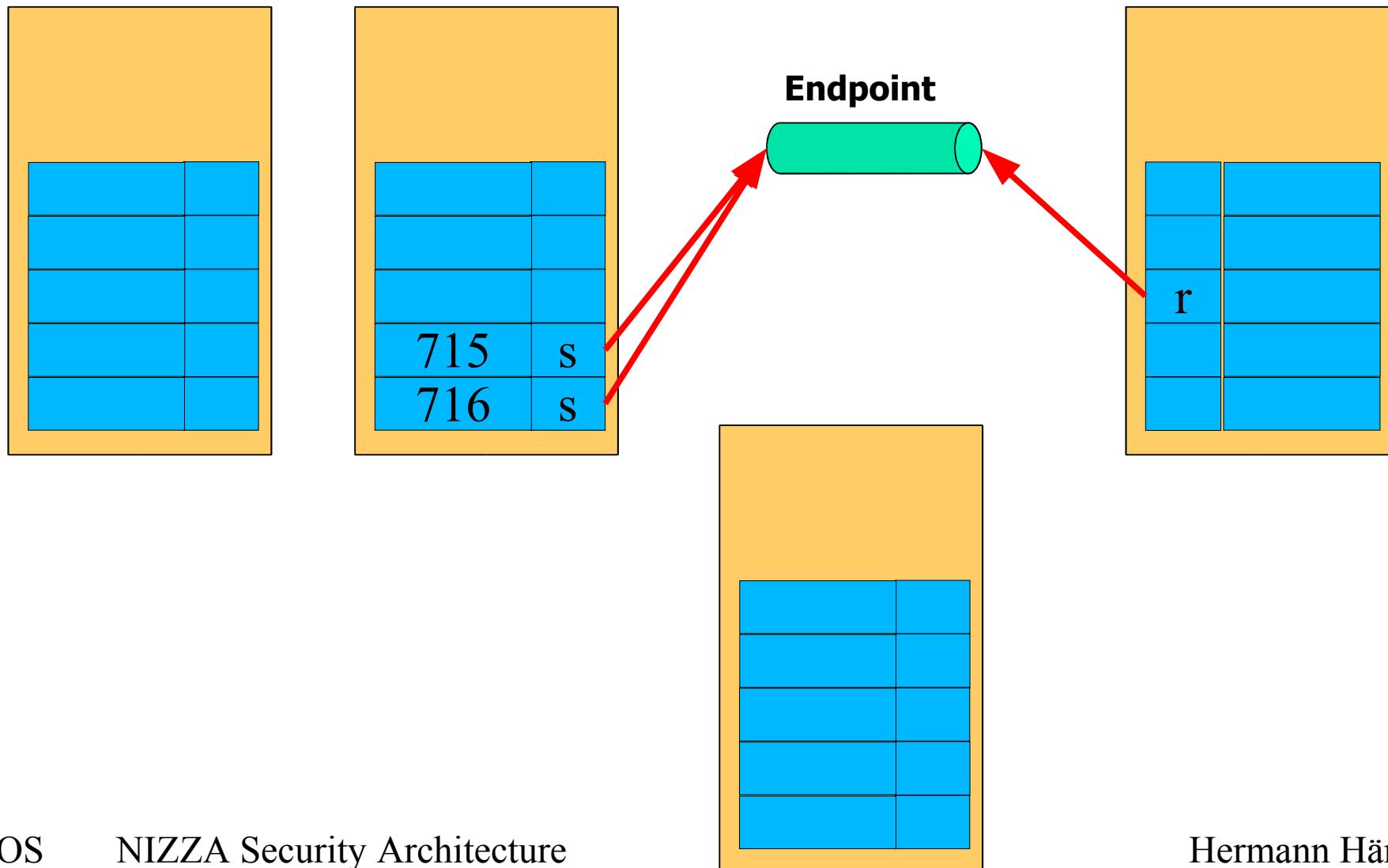


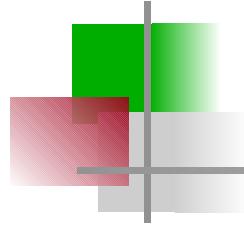
# Communication



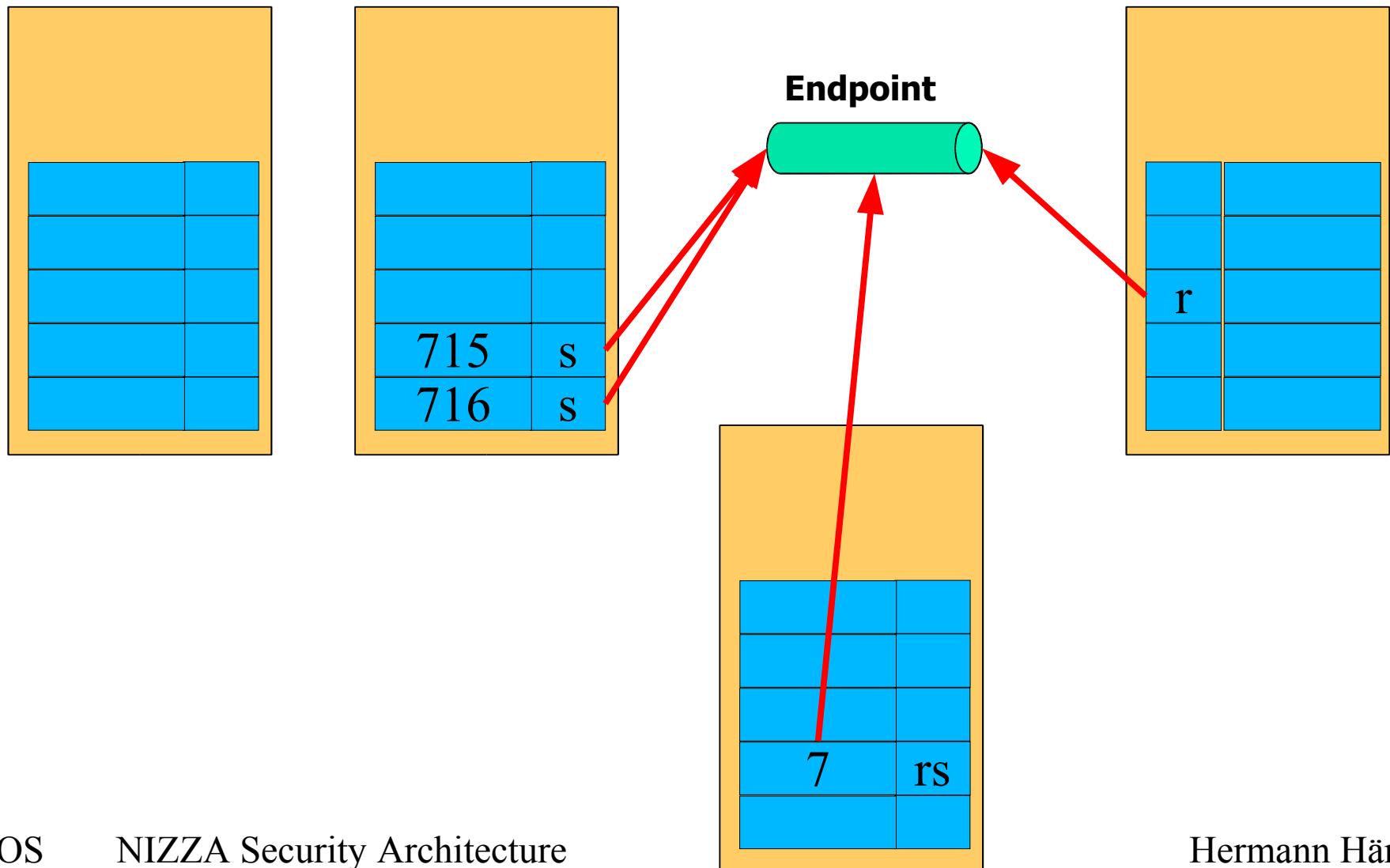


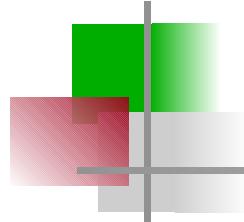
# Communication



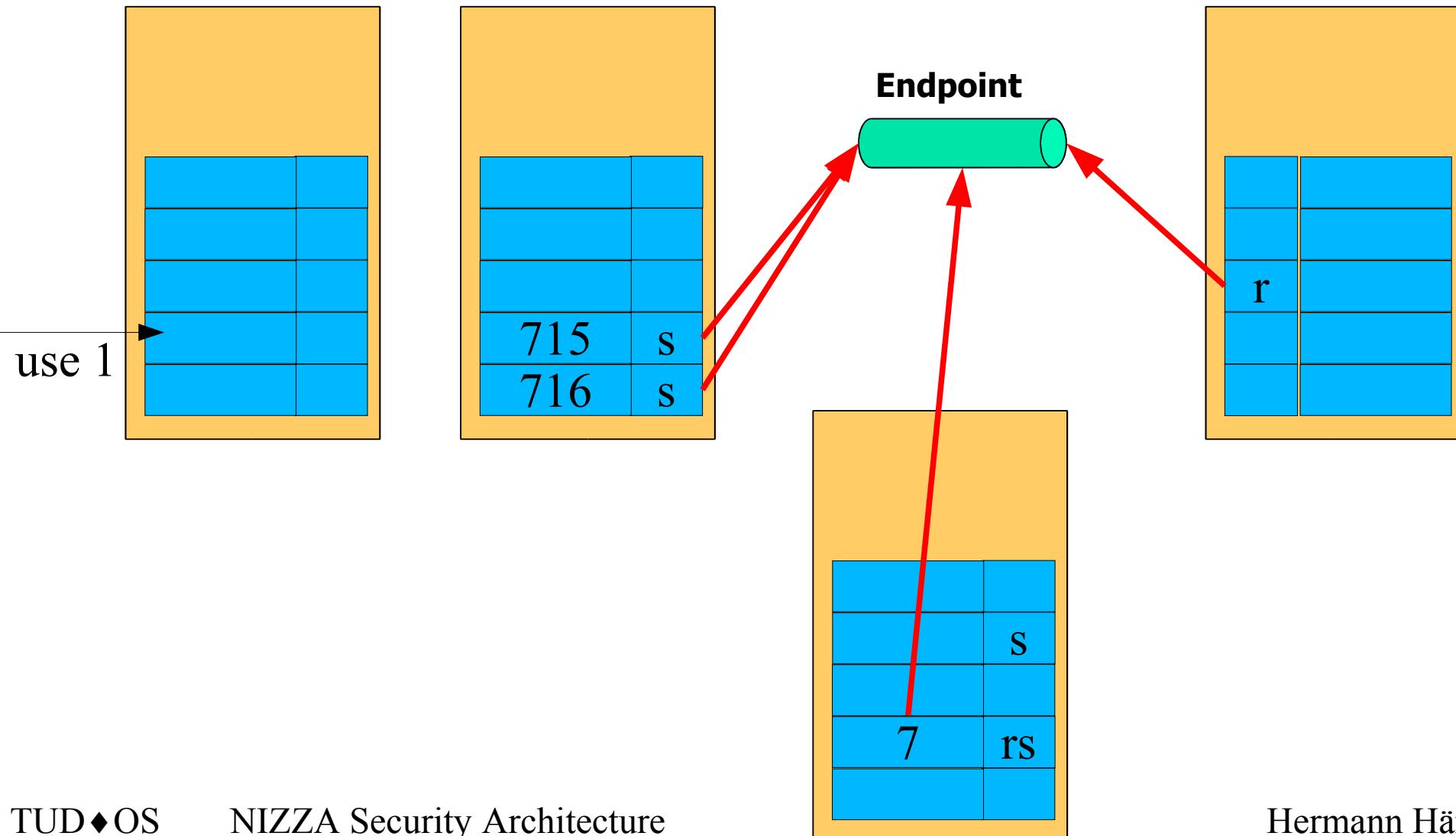


# Communication

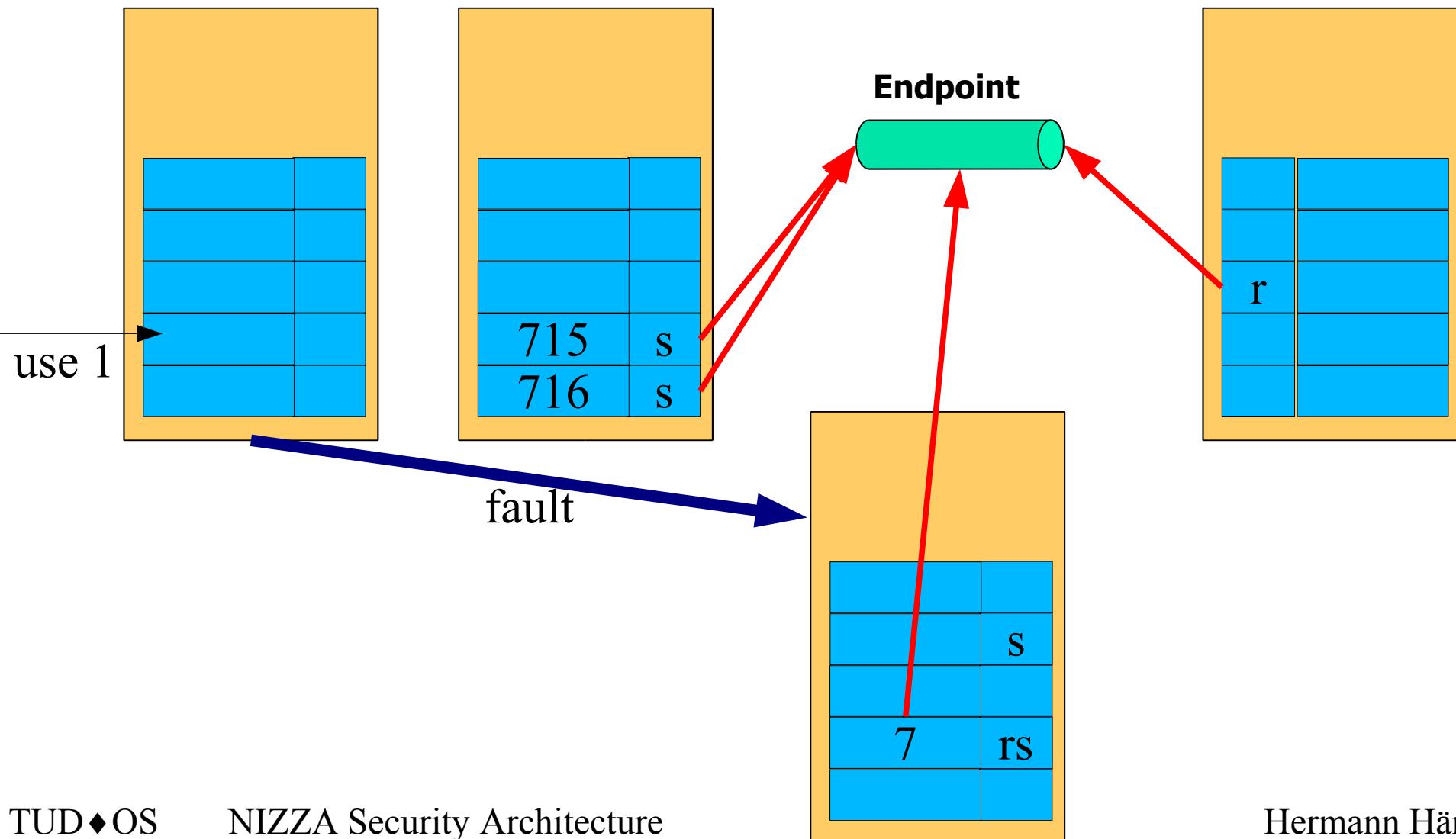




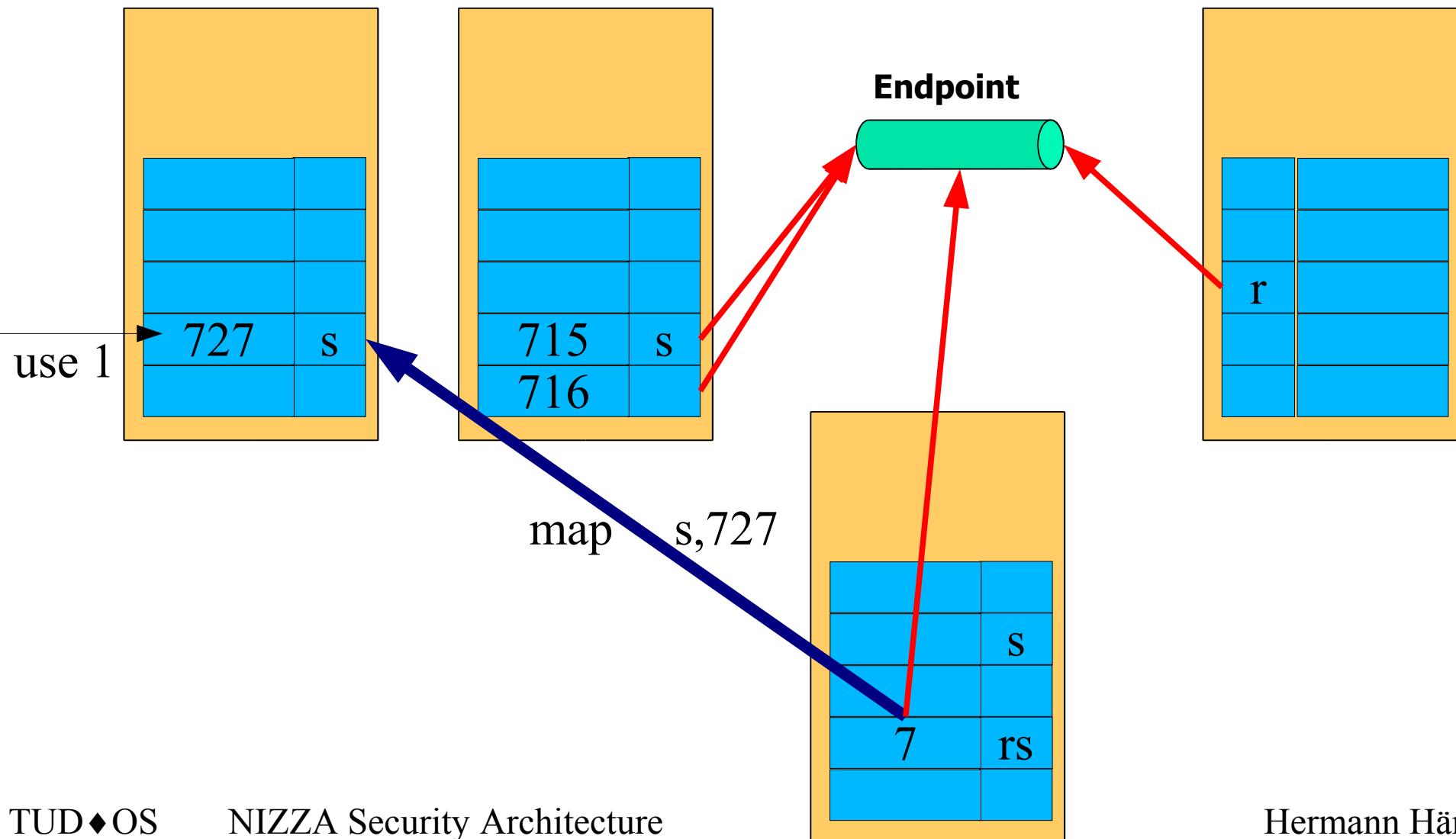
# Communication



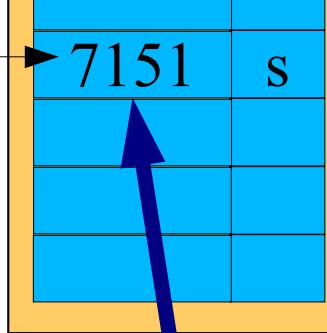
# Communication



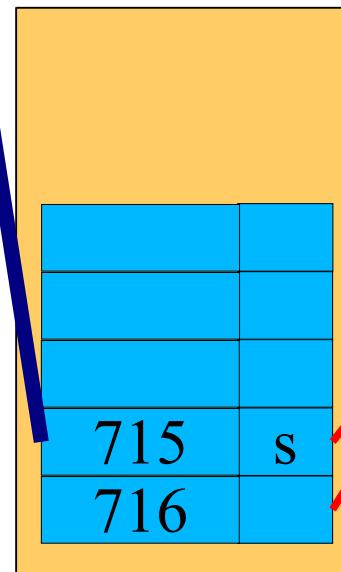
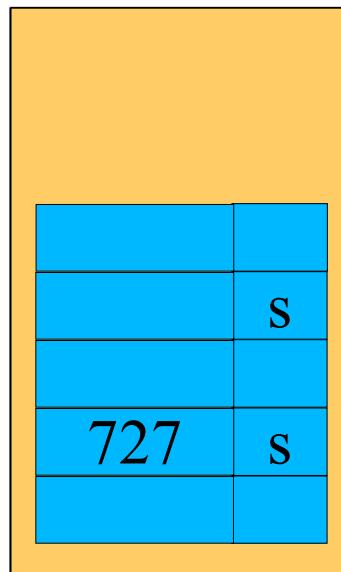
# Communication



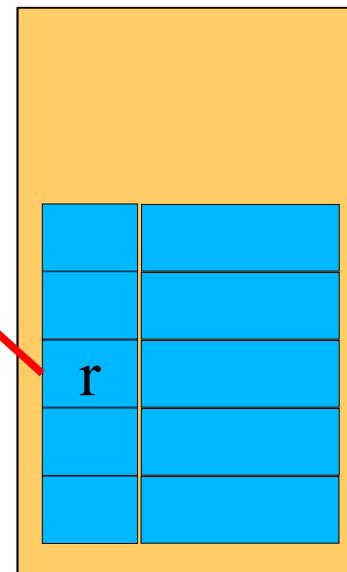
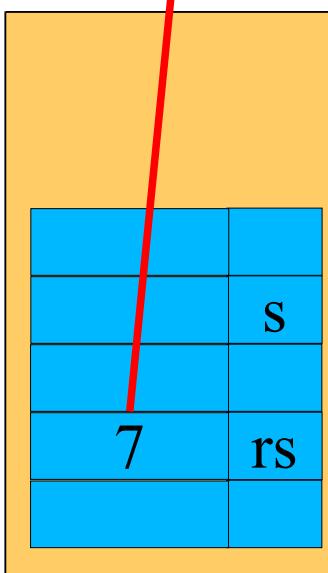
use 3



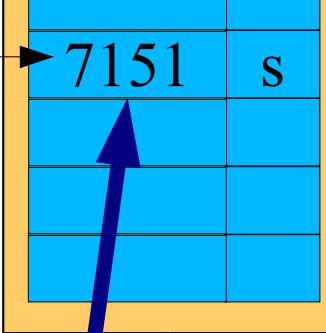
map s, 7151



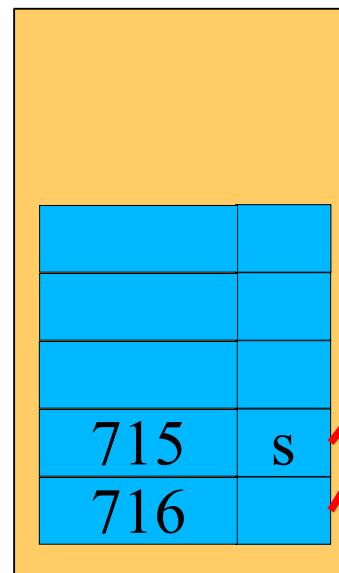
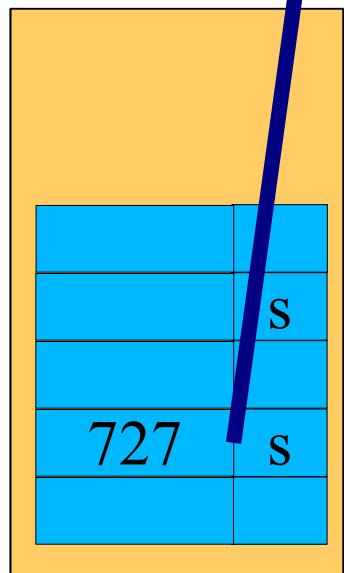
**Endpoint**



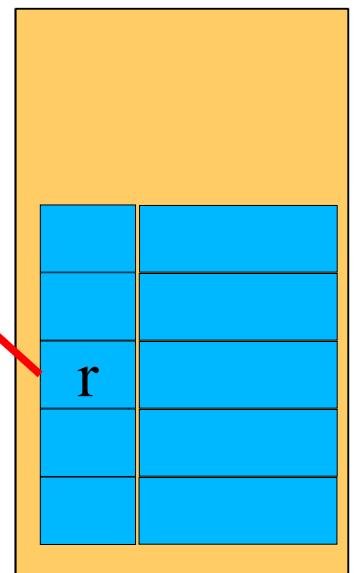
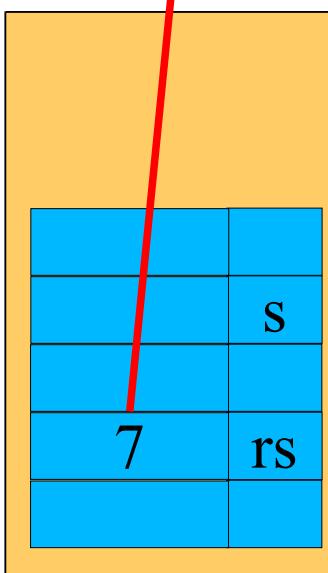
use 3



map s, 7271



Endpoint



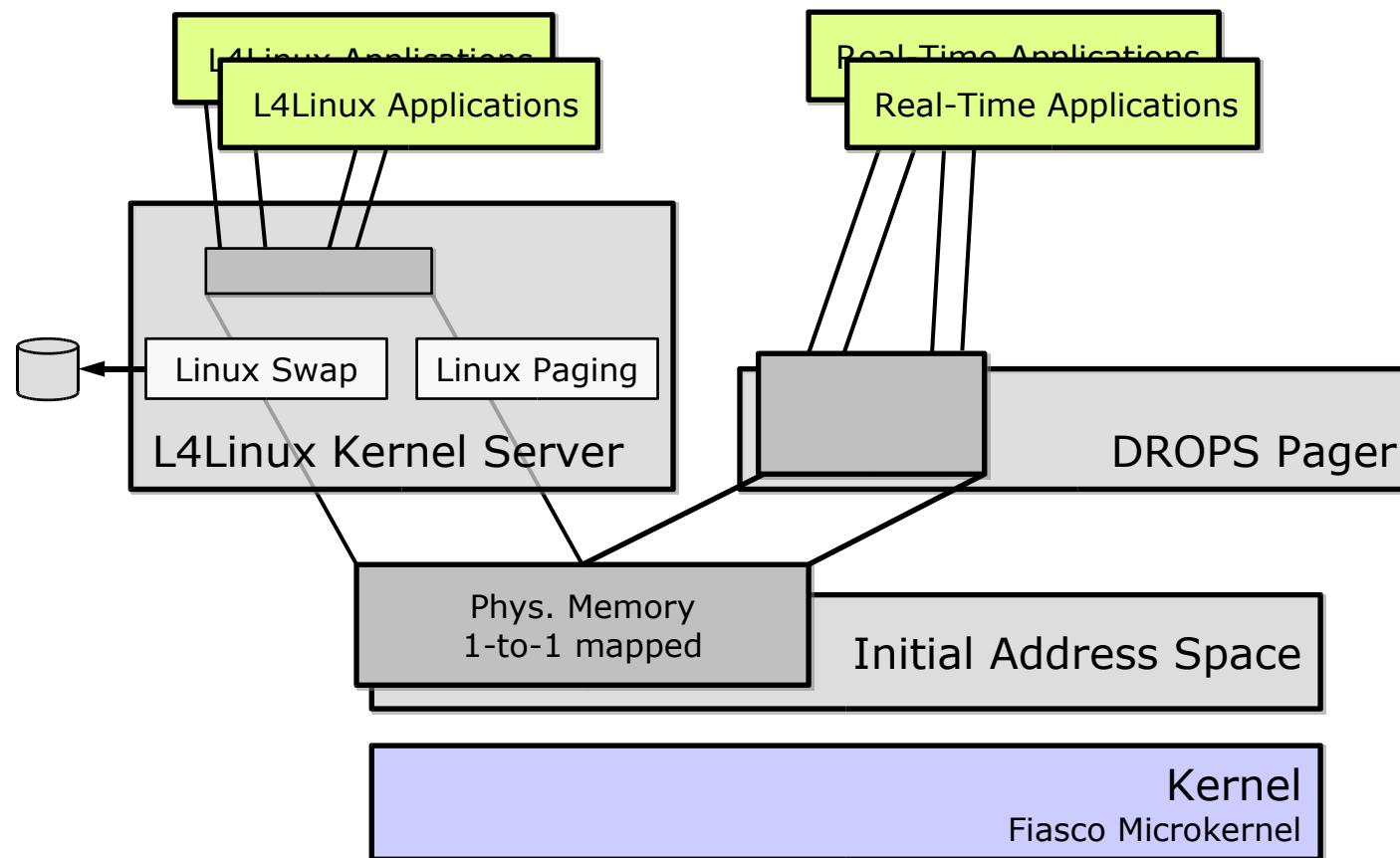


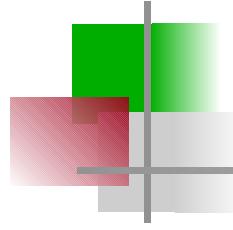
Linux Apps

L<sup>4</sup>Linux

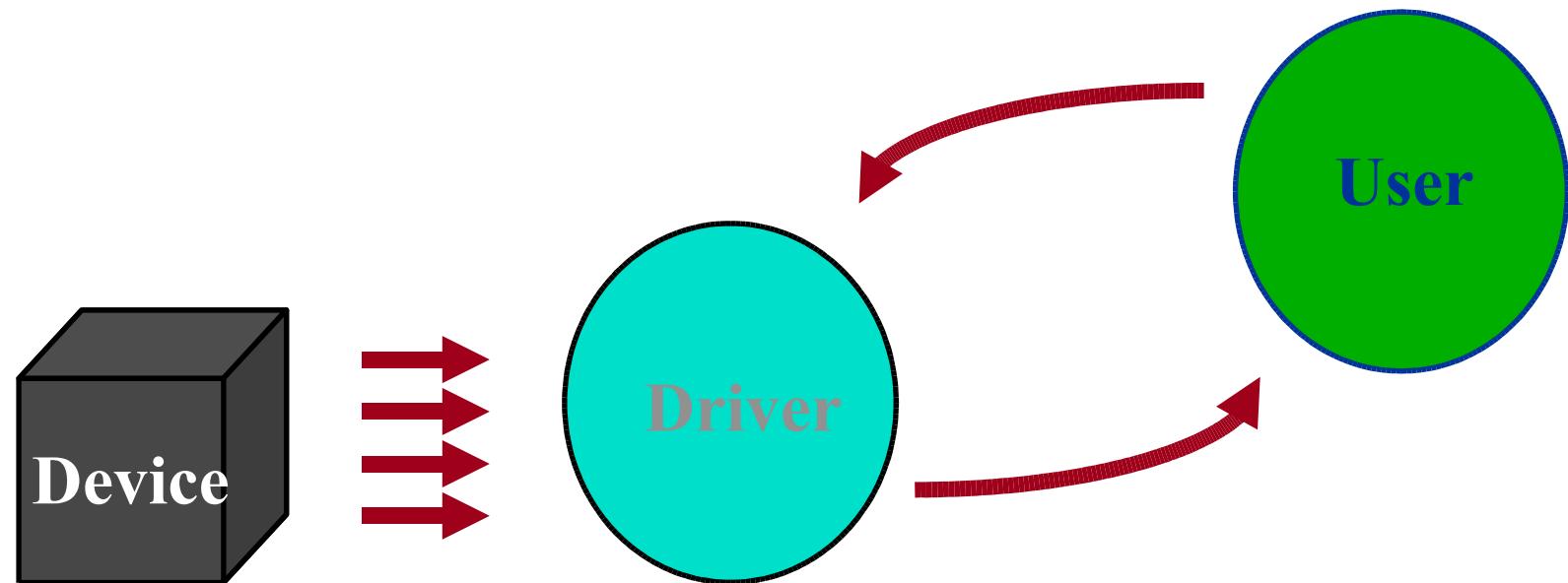
L4 Fiasco

# Pagers





# Linux “top halves” as threads

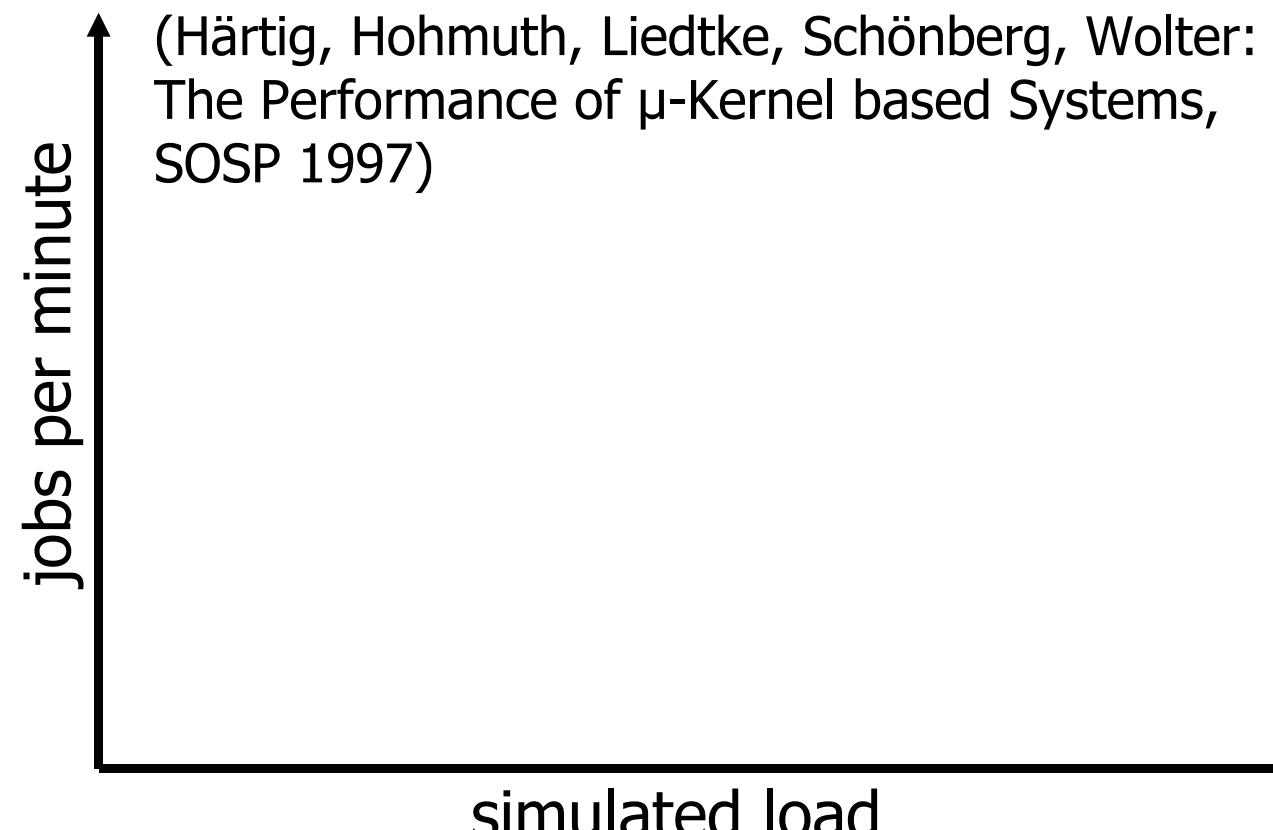


# Performance

Time-Sharing  
Applications

L<sup>4</sup>Linux

Fiasco

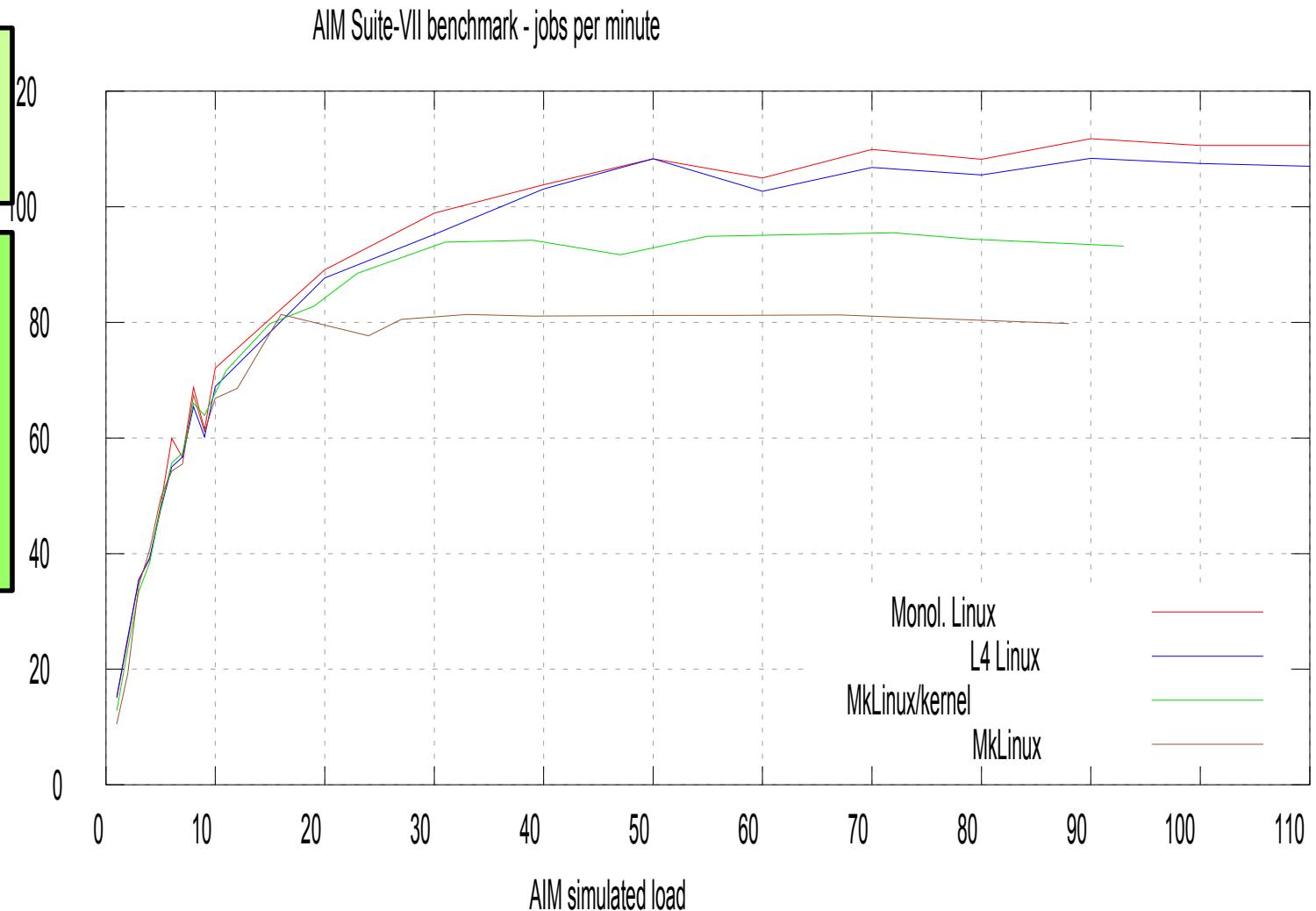


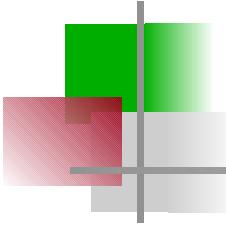
# L<sup>4</sup>Linux compared to MACH

Time-Sharing  
Applications

L<sup>4</sup>Linux

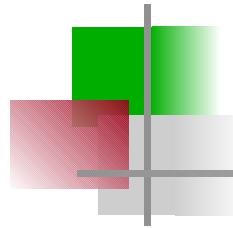
Fiasco



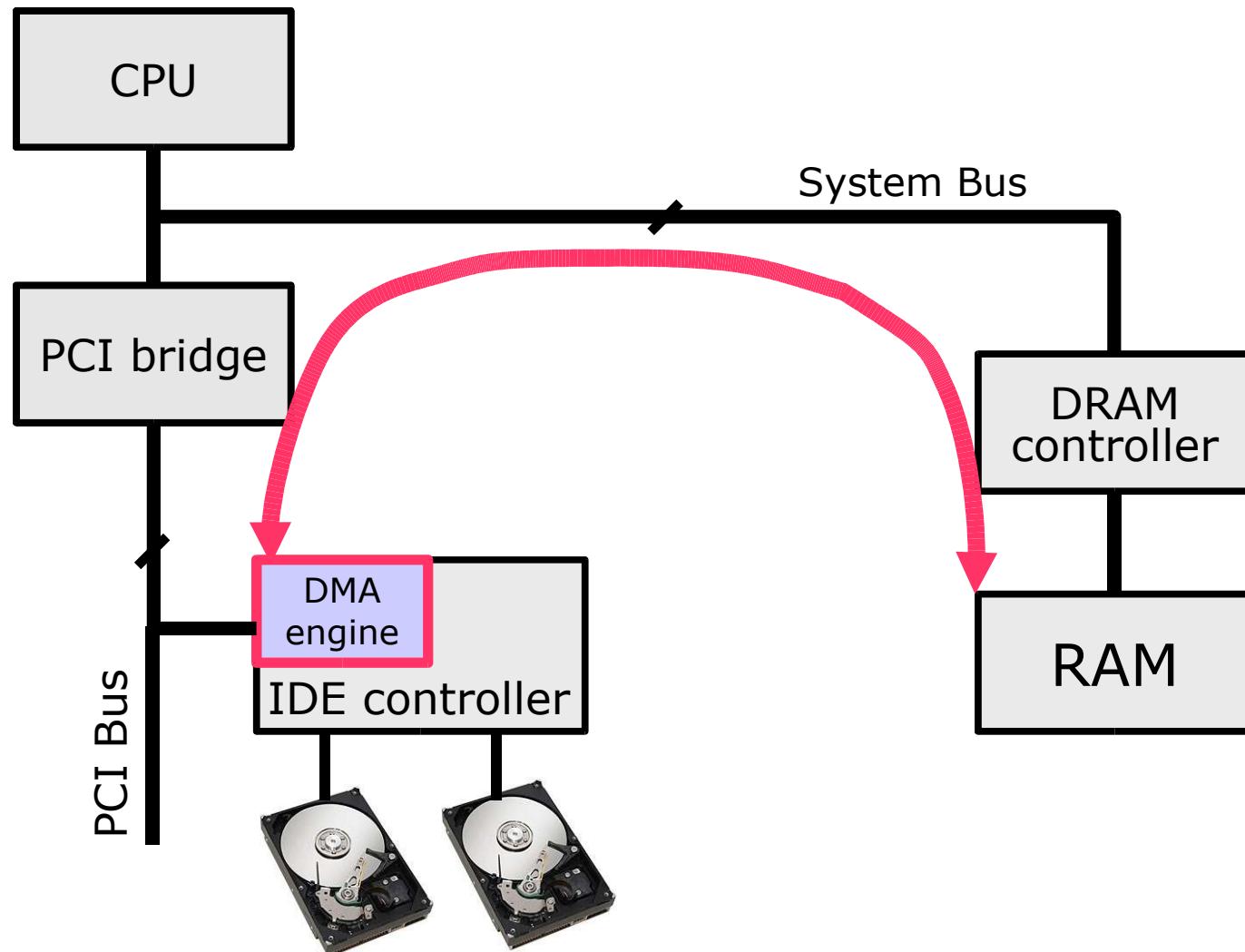


# The DMA Problem

- separation is enforced by MMU
  - devices access memory using bus master DMA
  - DMA uses physical addresses  
(on most architectures)
- malicious devices  
(or malicious device drivers, firmware!)  
can access/modify all components of a system



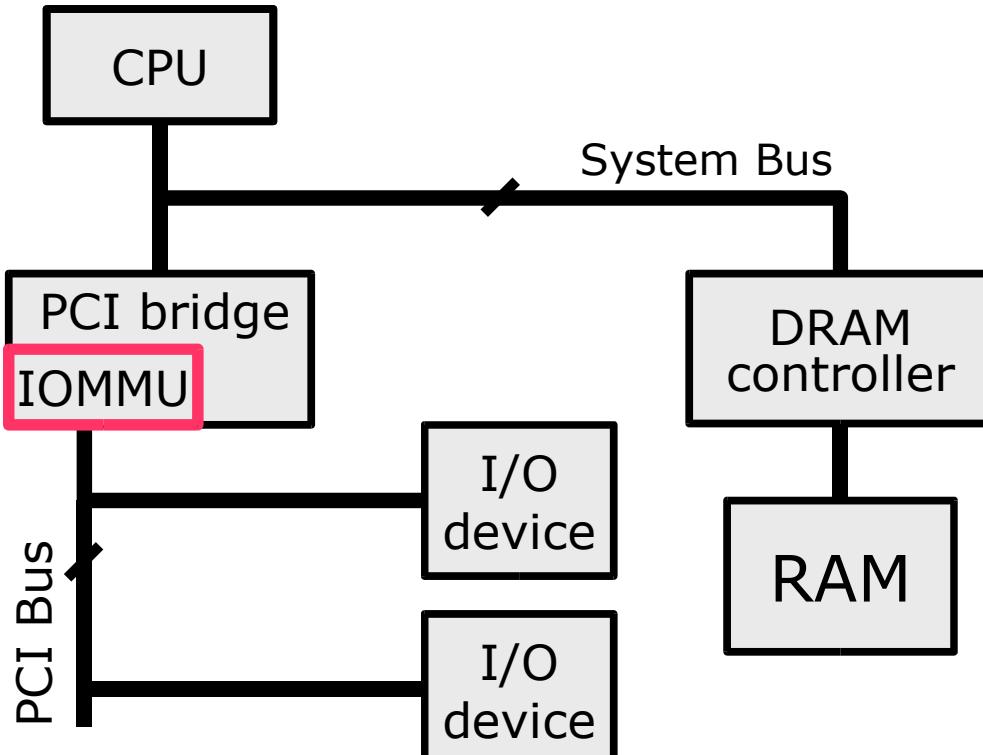
# The DMA Problem



# IOMMU

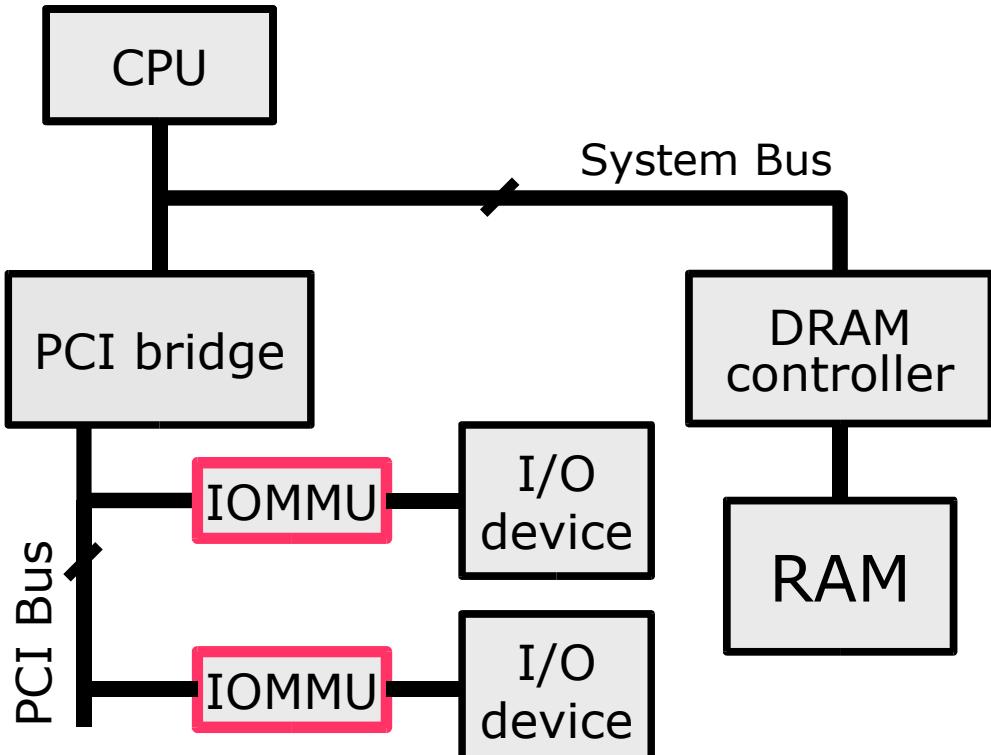
## per-bus IOMMU:

- IA64 chipsets
- Opteron



## per-device IOMMU:

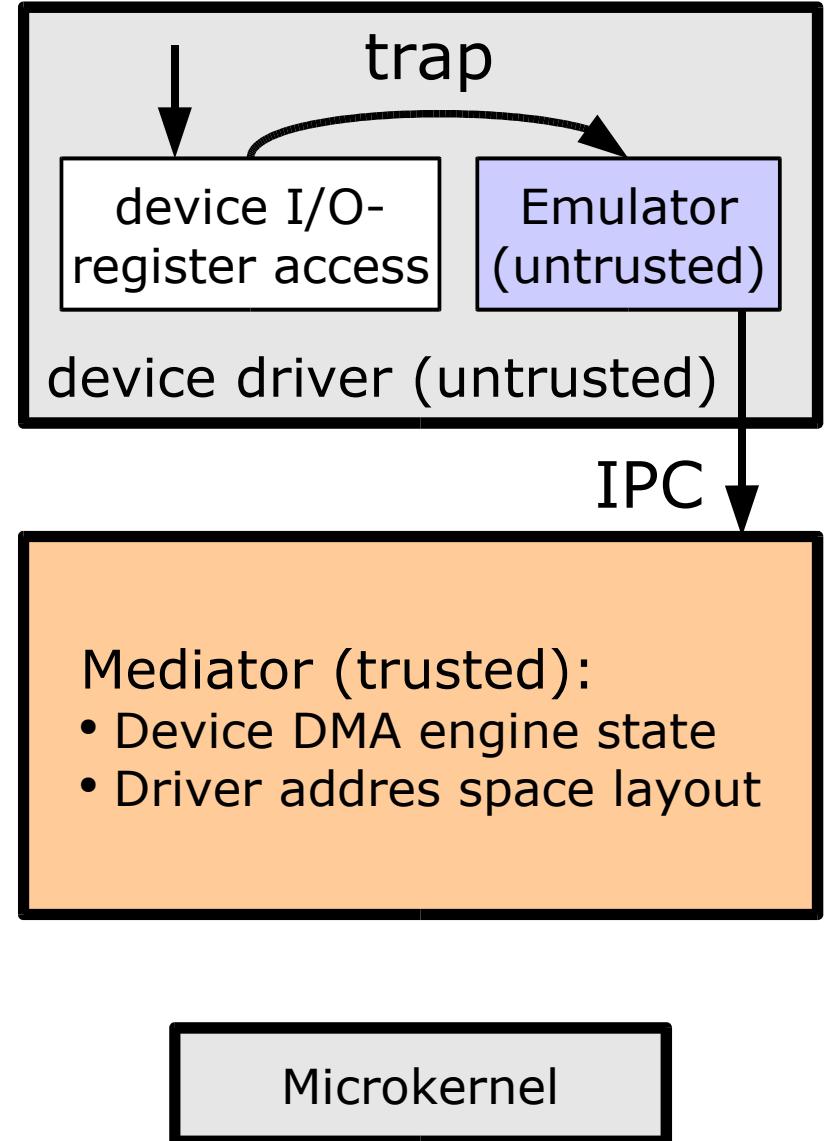
- **n/a** in current hardware!

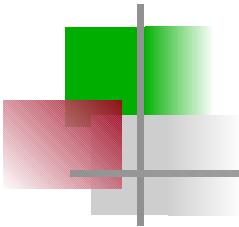


# IOMMU in Software

Make sure that device cannot perform malicious DMA:

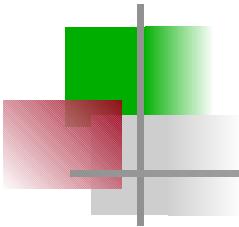
- **Trap** read/write access to I/O-registers of the device
- An **emulator** (untrusted) determines size and value to be read/written
- A **mediator** (trusted) checks and performs the access





# IOMMU in Software

- **Emulator:**
  - in driver's address space → untrusted
  - malfunction does only decrease availability of device
  - ~ 500 LoC
- **Mediator:**
  - trusted, own address space
  - specific for a device or a class of devices
  - ~ 300 LoC

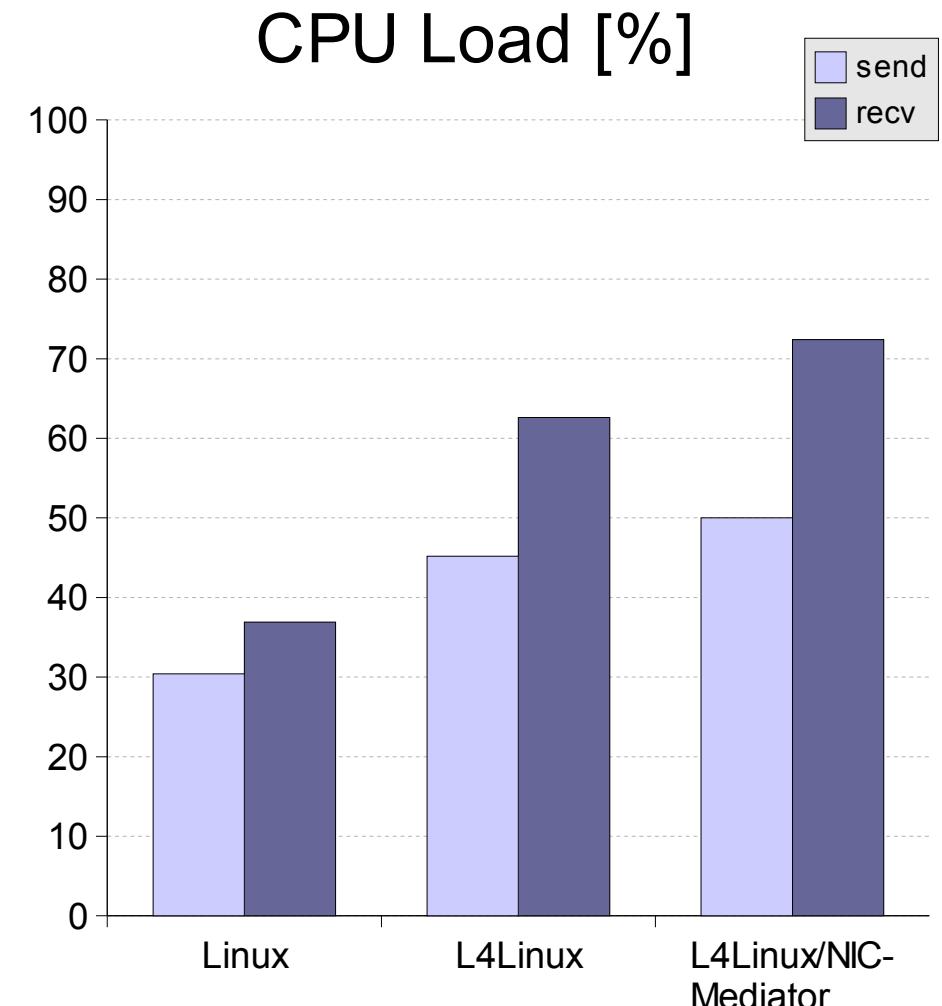
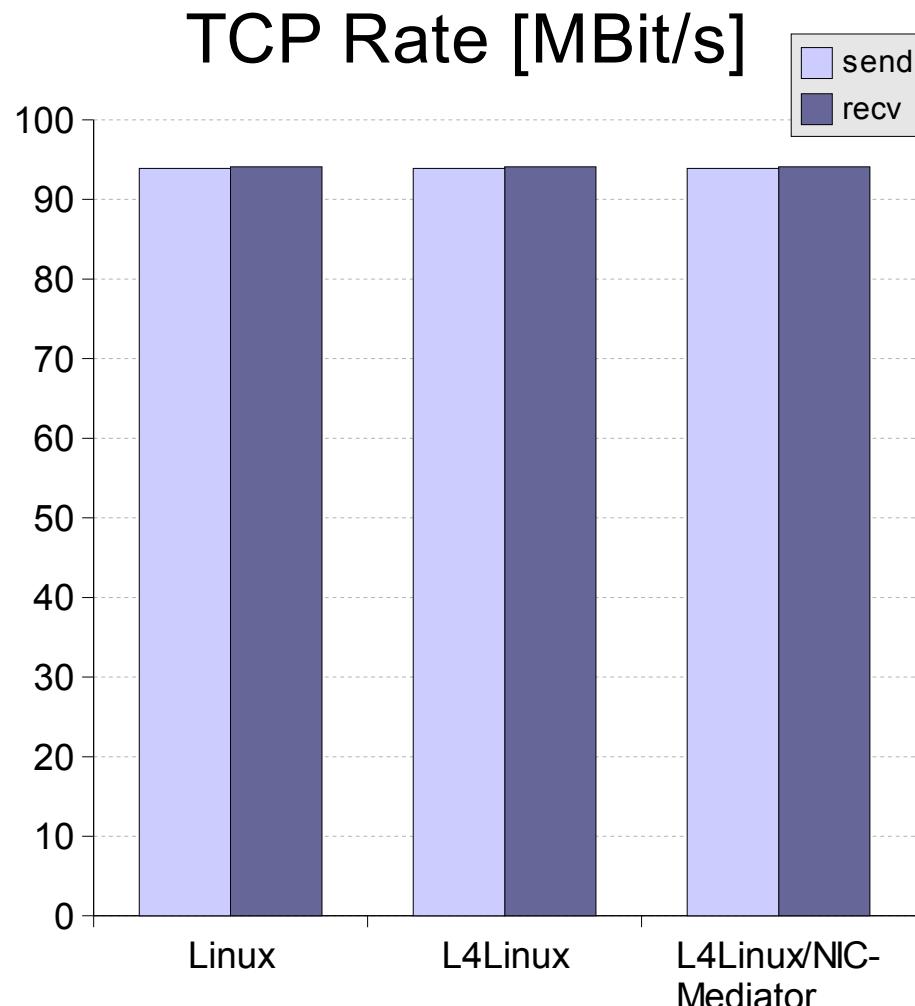


# IOMMU in Software

- Implemented for:
  - Fast Ethernet card (DEC Tulip 21143)
  - ATA Controllers
- Does not work for firmware-programmable devices:
  - private interface between device driver and firmware
  - no protection mechanism between firmware and device

# IOMMU in Software: NIC Performance

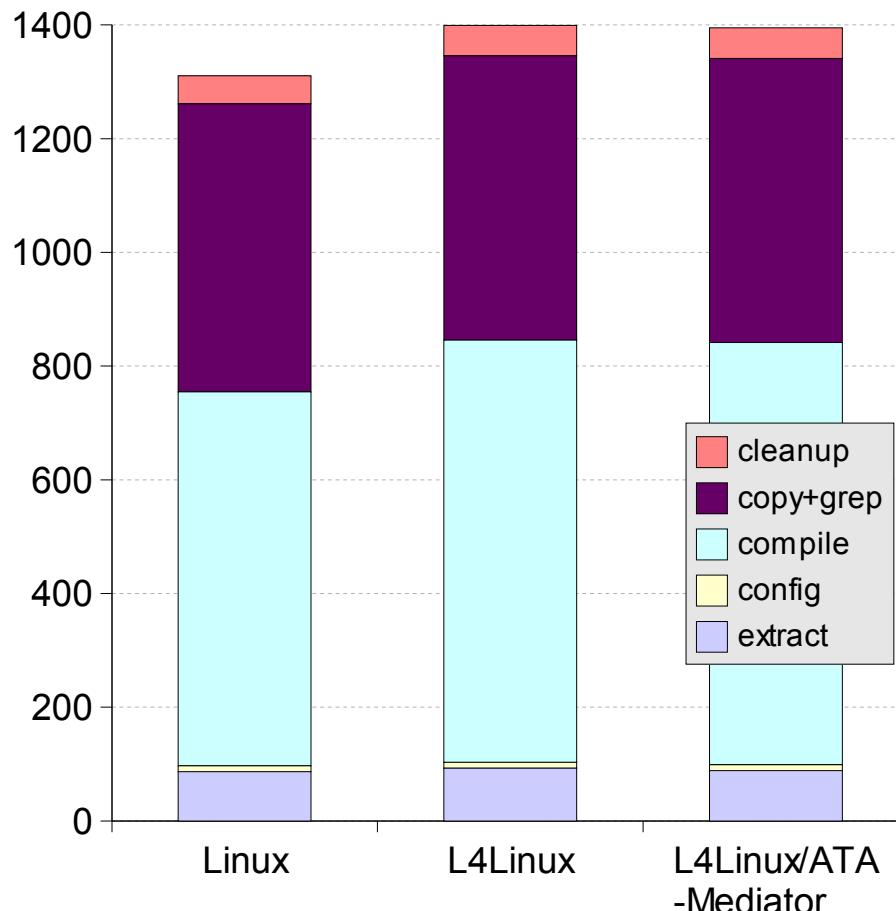
Pentium-III 800 MHz, Fast Ethernet



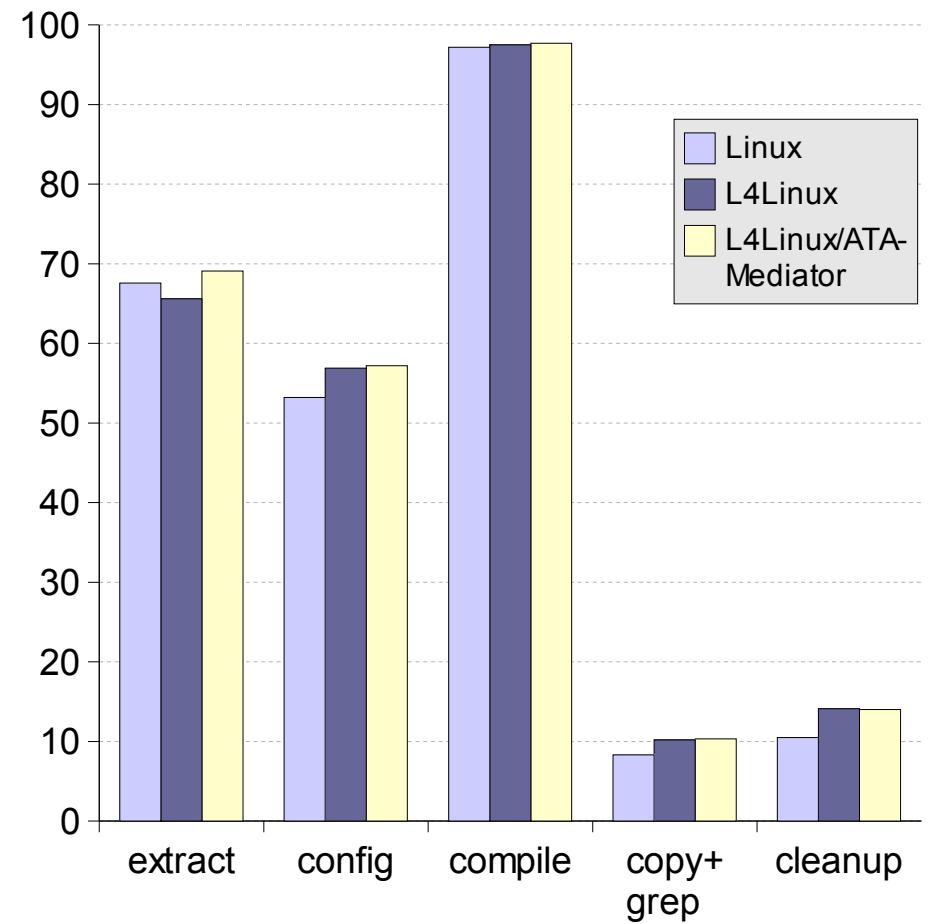
# IOMMU in Software: ATA Performance

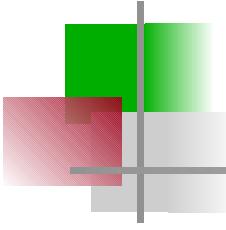
Pentium-III 800 MHz, VIA82C586 ATA Controller

Total Time [s]

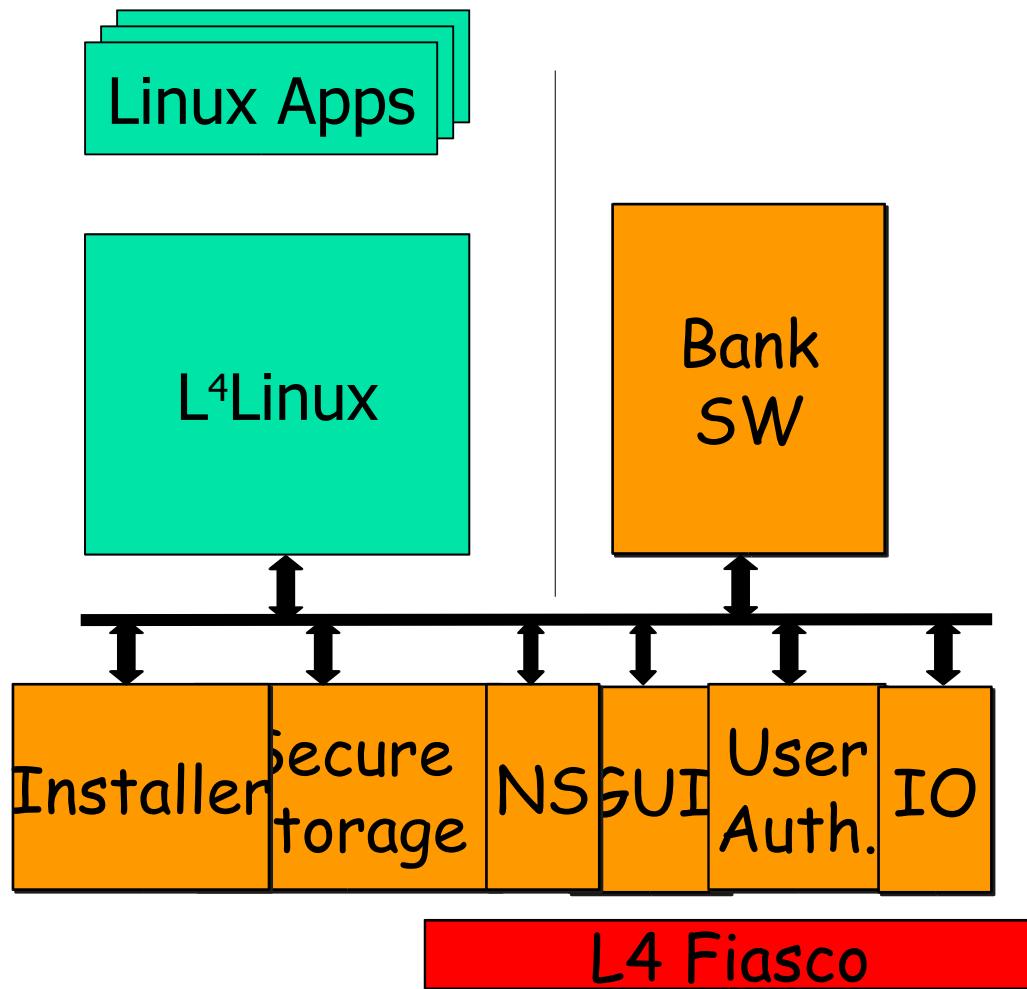


CPU Load [%]





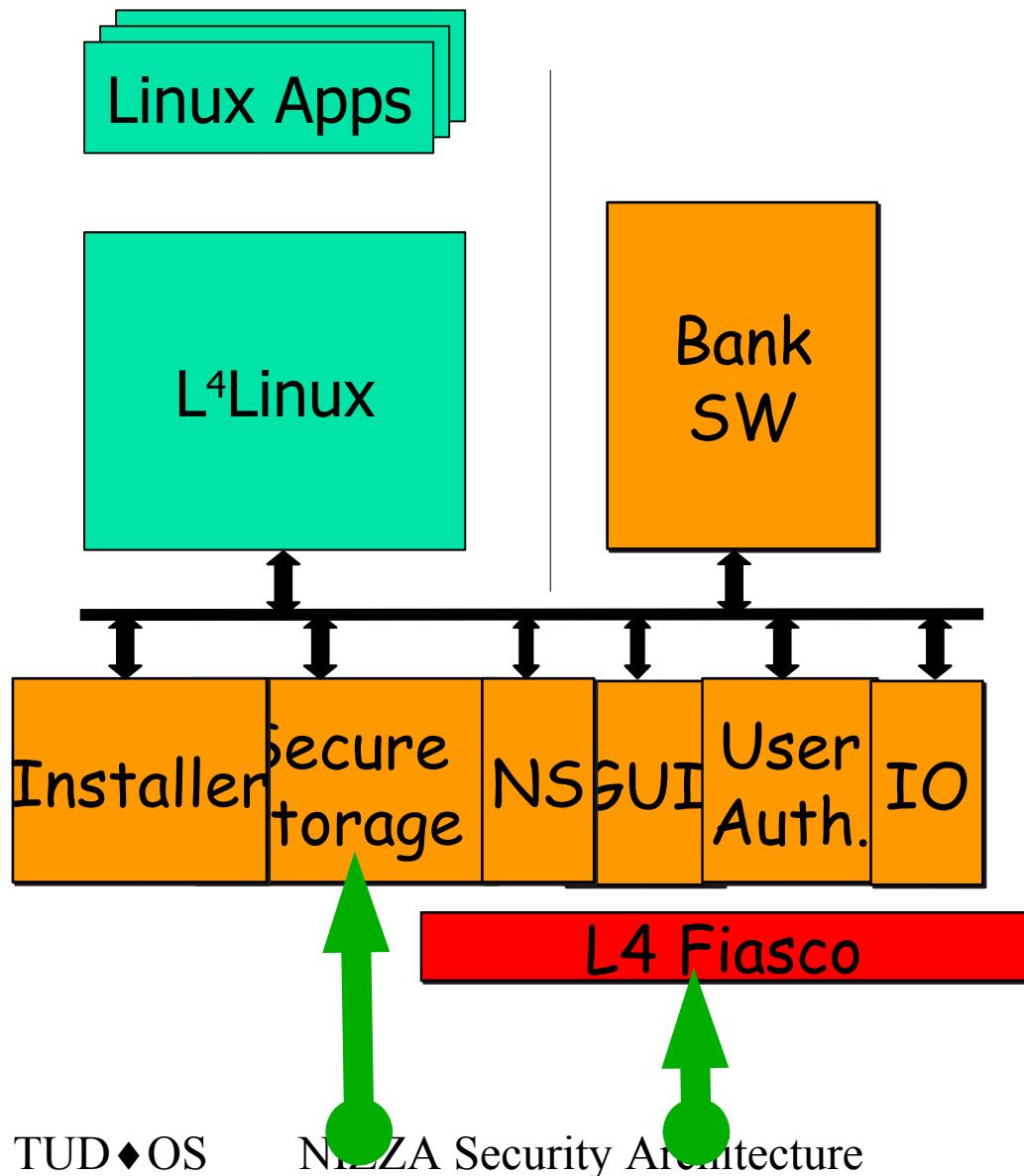
# Secure Booting, Remote Attestation and Trusted Path



Authentic application/system:

- how does the remote bank know ?
- how does the local client/user know ?

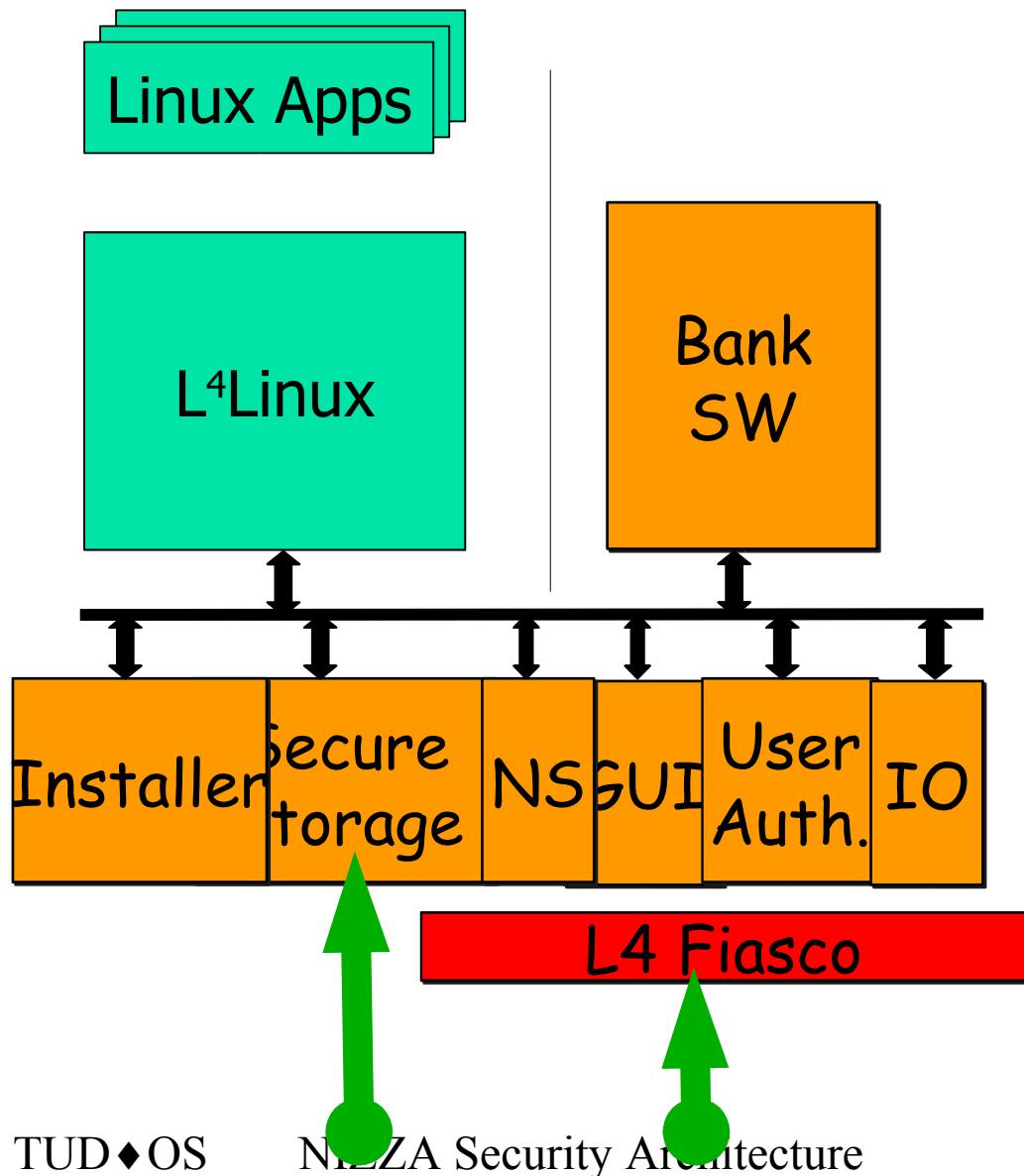
# Secure Booting, Remote Attestation and Trusted Path



Authentic application/system:

- how does the remote bank know ?
- how does the local client/user know ?

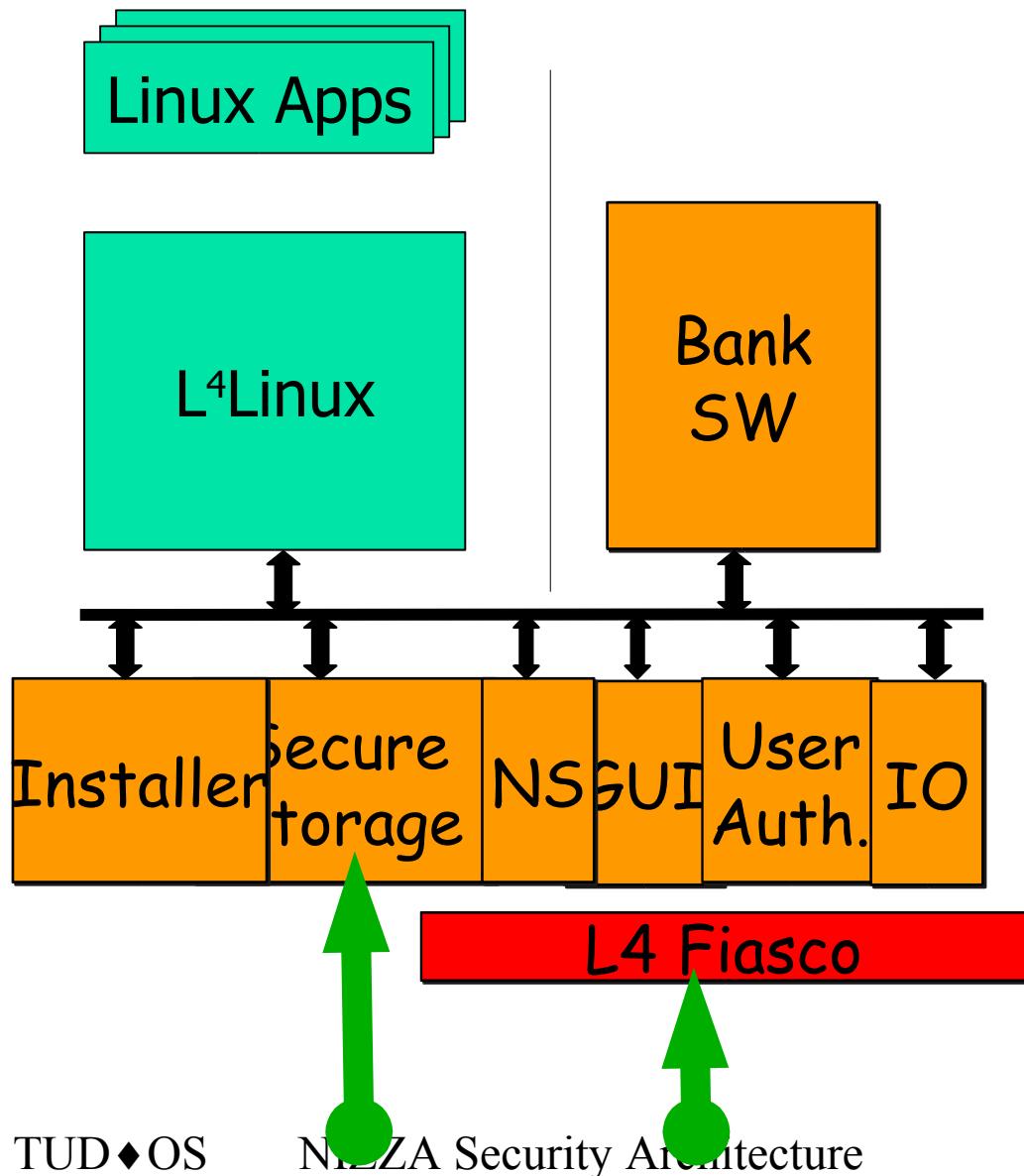
# Secure Booting, Remote Attestation and Trusted Path



Authentic application/system:

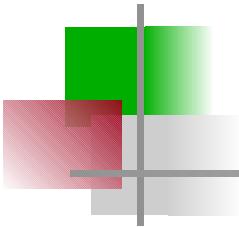
- how does the remote bank know ?
  - attestation protocol up to Nizza trusted platform
  - mediate other communication thru trusted installer

# Secure Booting, Remote Attestation and Trusted Path



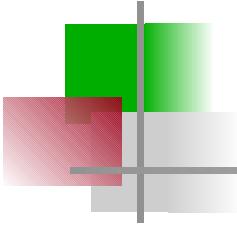
Authentic application/system:

- how does the local client/user know ?
  - attestation protocol up to Nizza trusted platform
  - indicate “red/green”
  - handover to DOpE



# Secure Storage with small TCB (future)

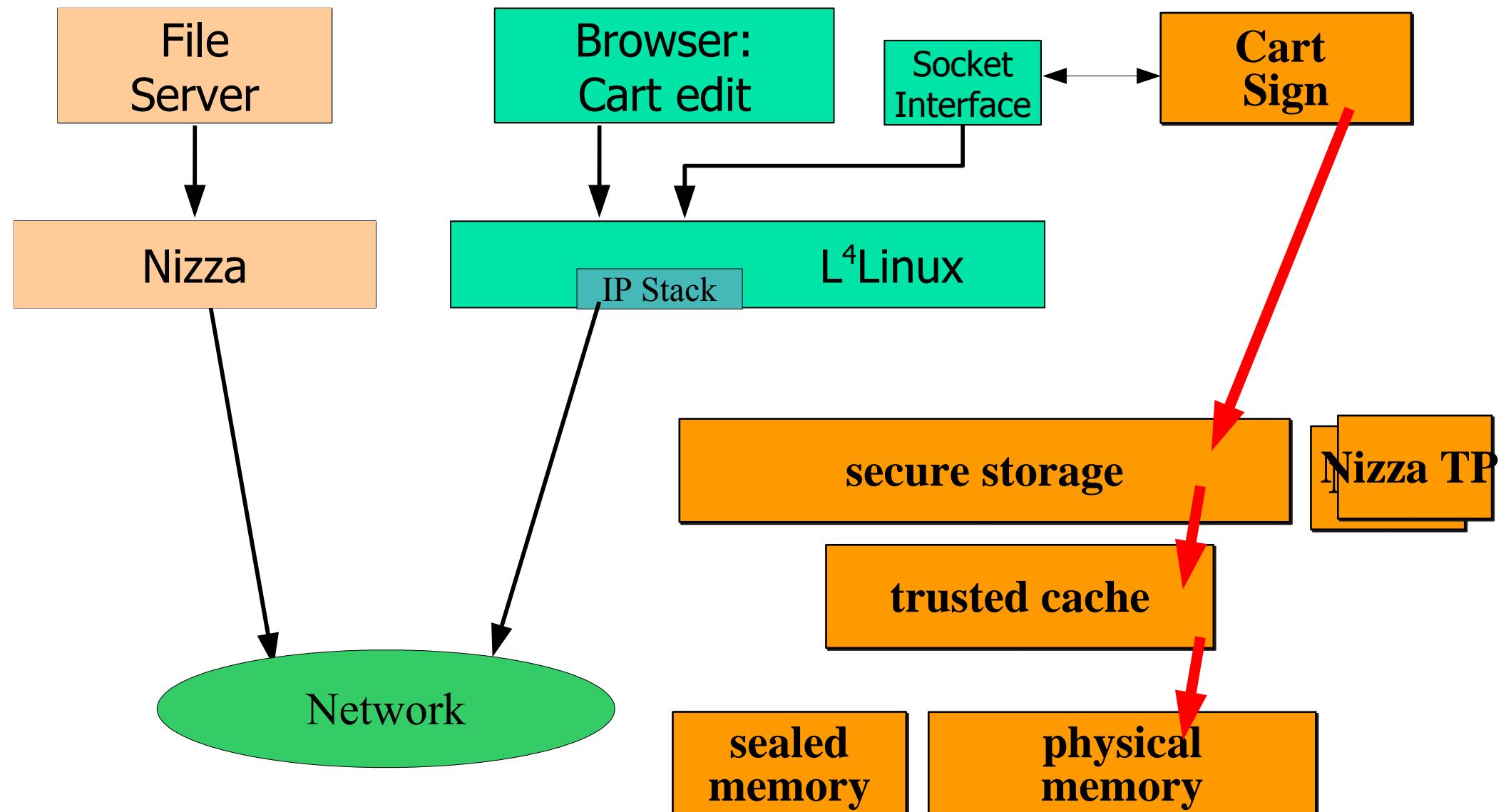
- objectives:
  - security
    - confidentiality, integrity,
    - recoverability
    - availability
  - system security
    - small TCB
    - attacks:
      - theft/loss of device
      - full penetration of L4Linux
-



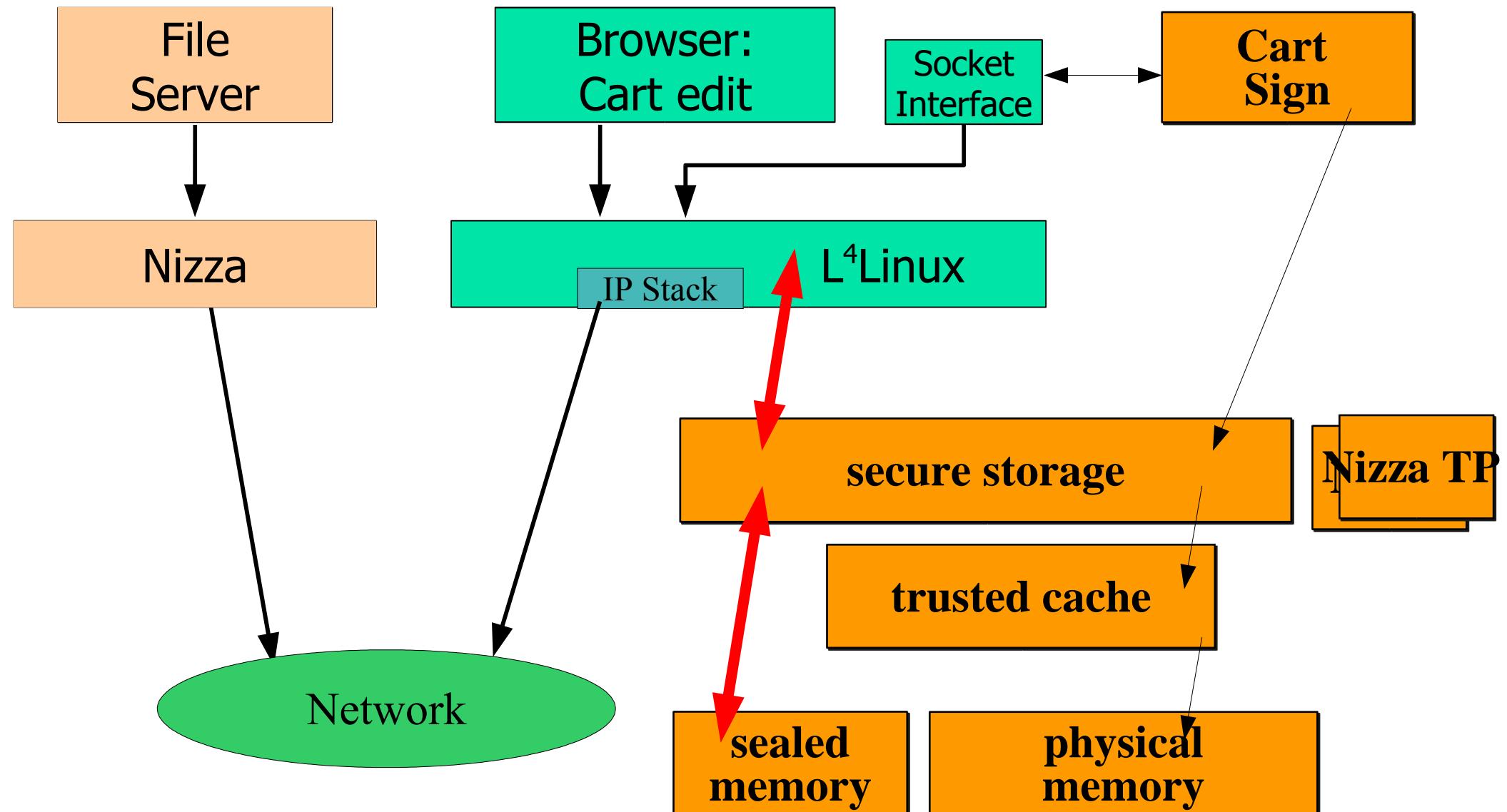
# Secure Storage with small TCB (future)

- techniques:
  - use Sealed Memory as key storage
  - reuse L4Linux file system as mass storage
  - use trusted file server for recoverability
  - use resource allocation for availability

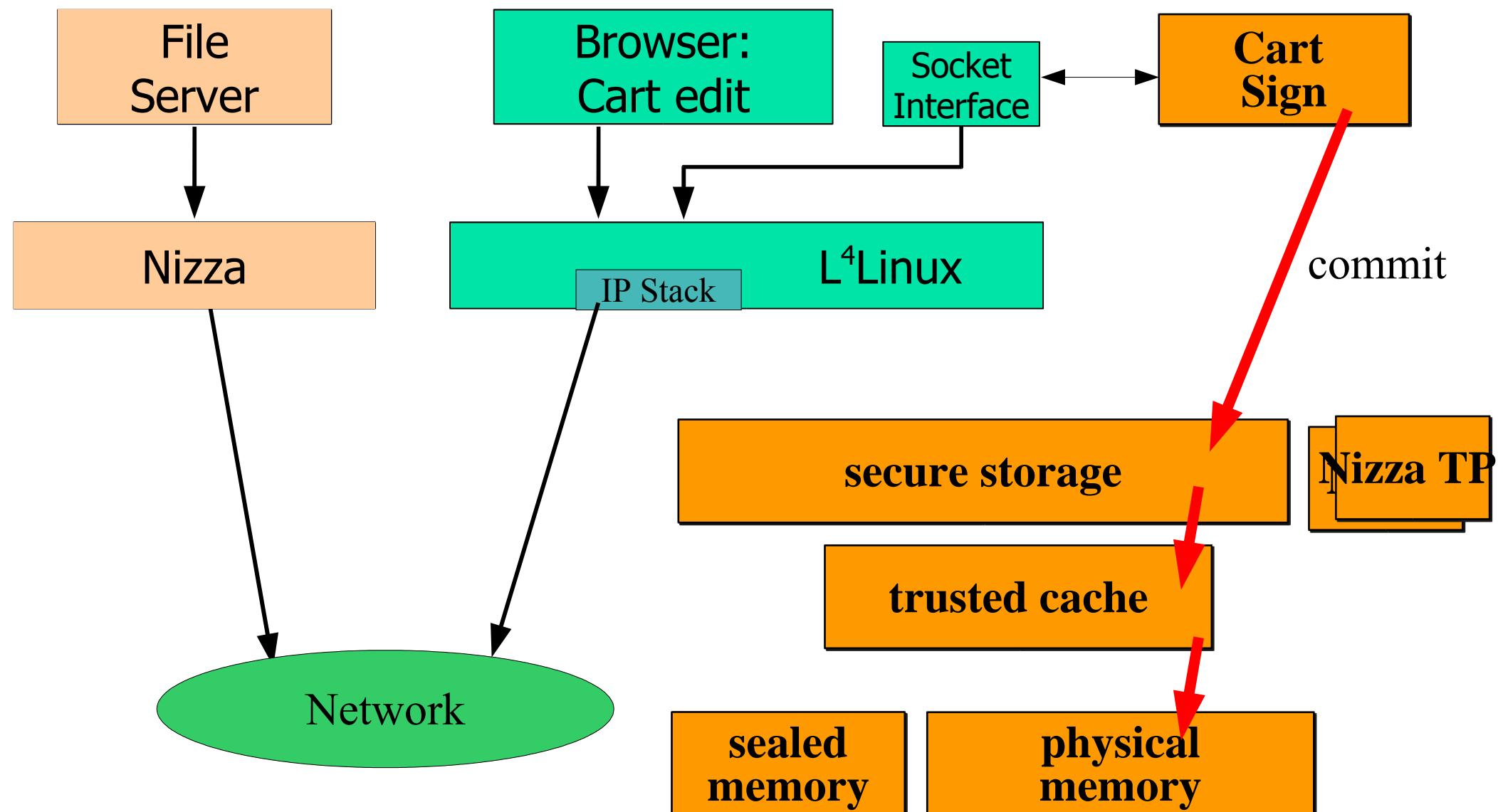
# Confidentiality and Integrity



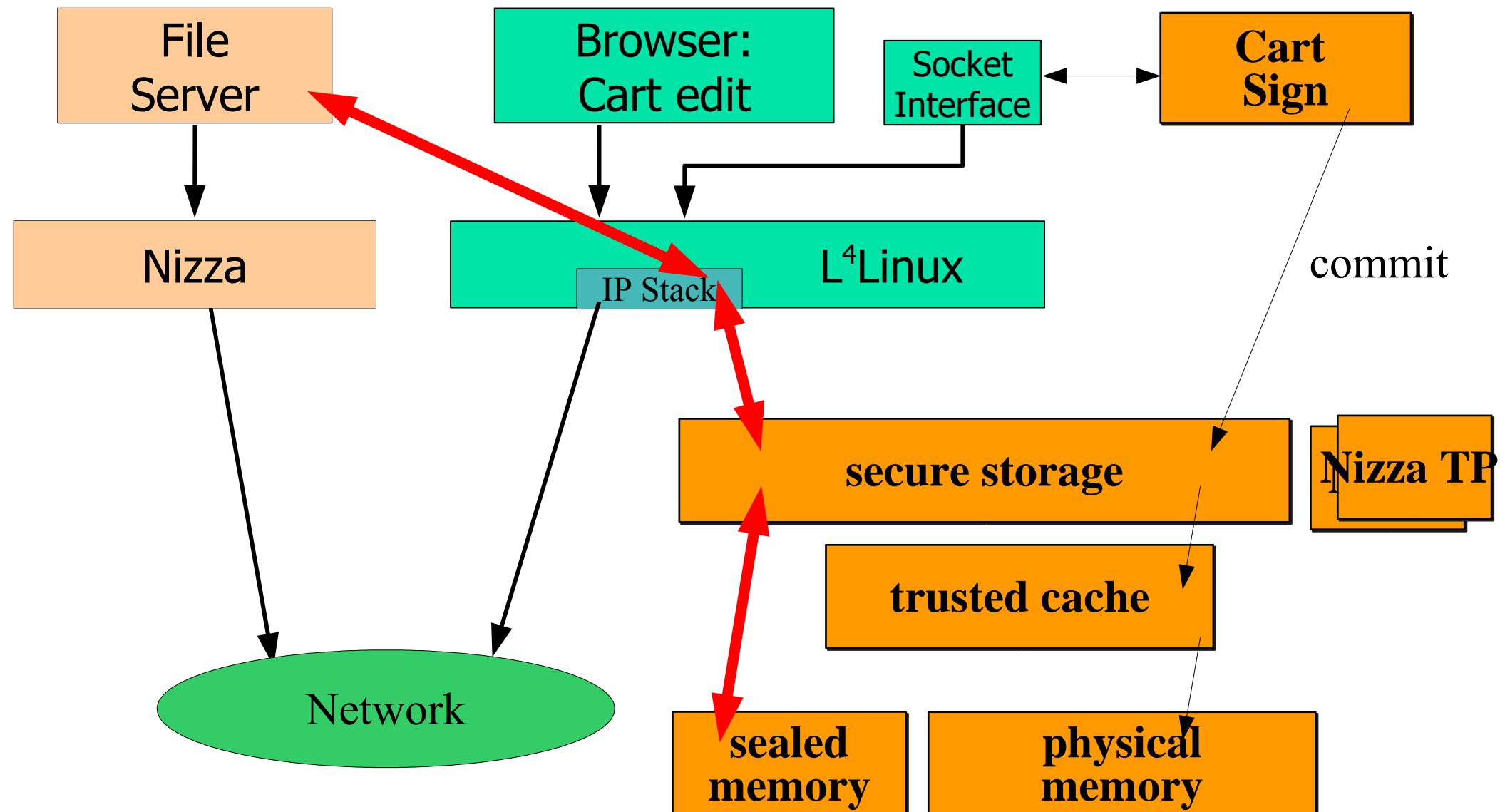
# Mass Storage (File System) not part of TCB



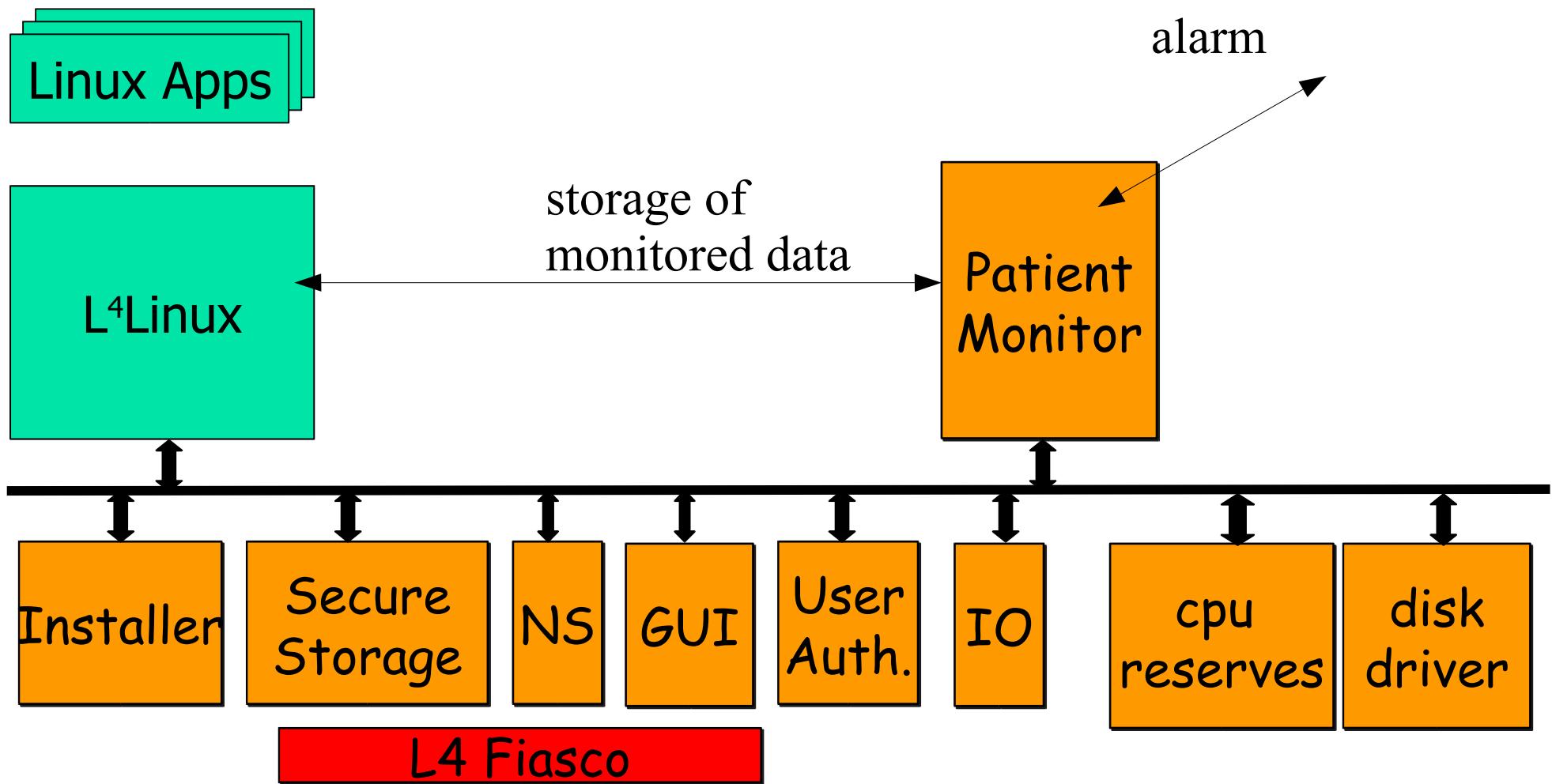
# Recoverability



# Recoverability



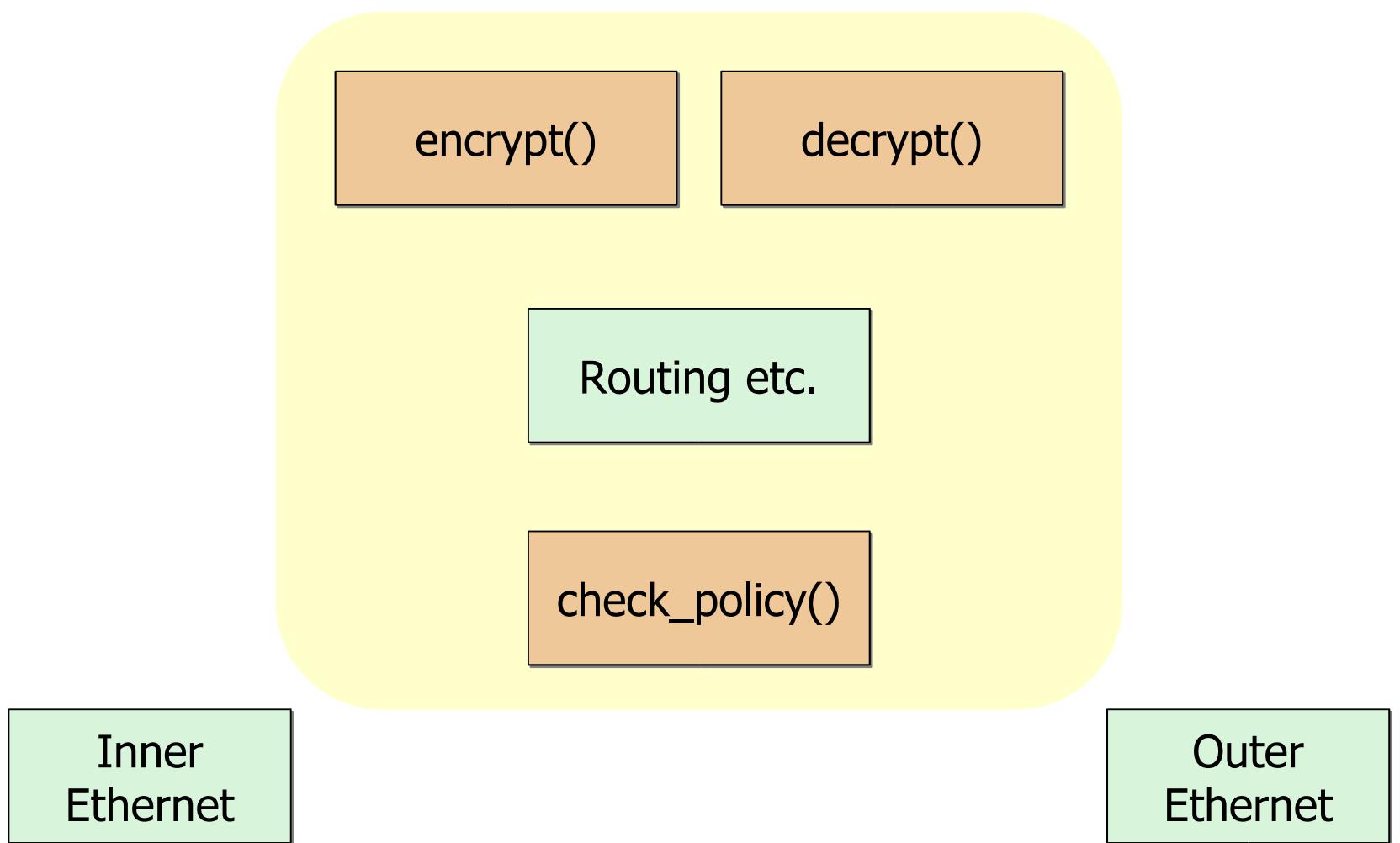
# Availability for Secure Storage



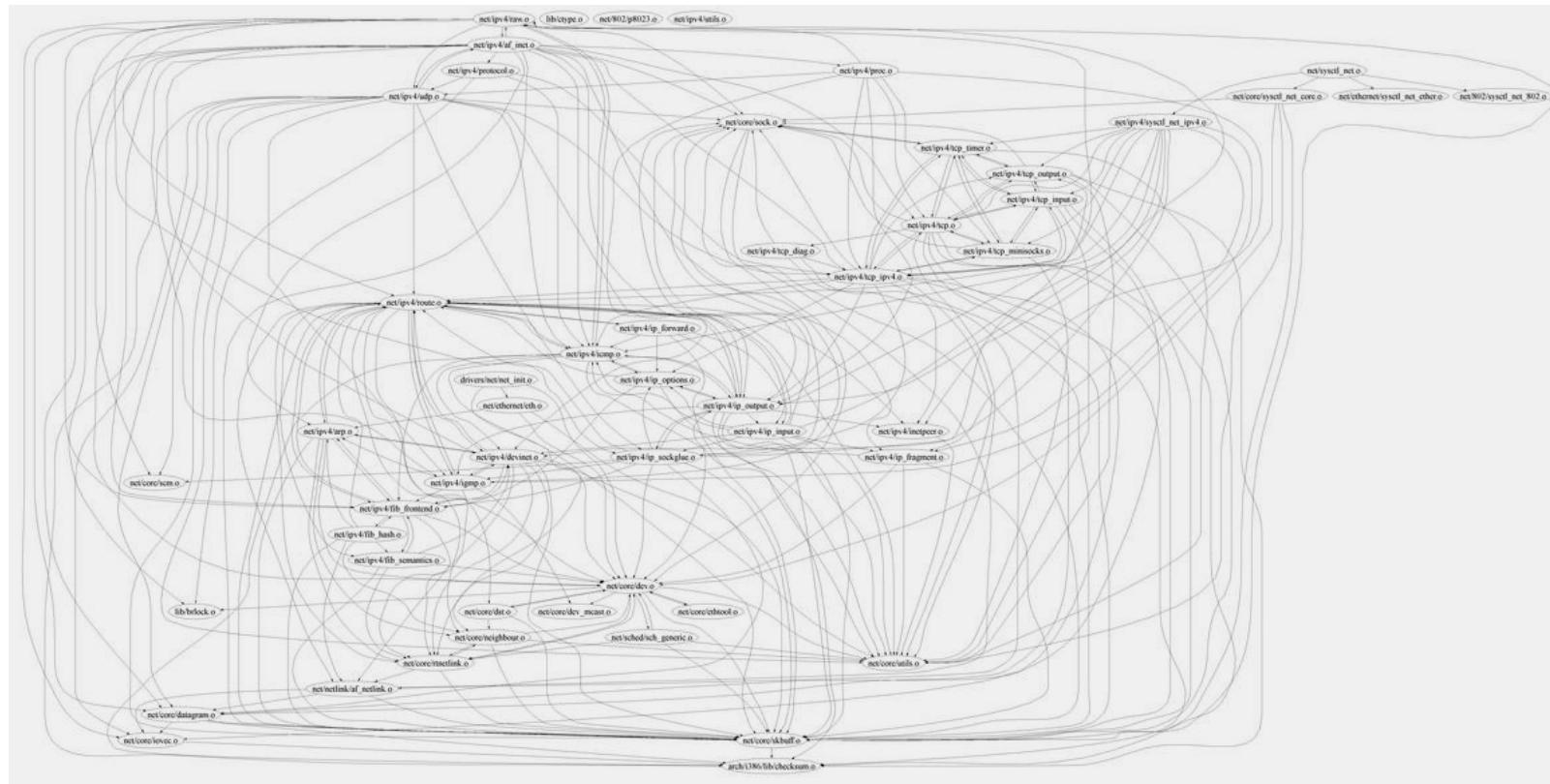
- VPN Box
- reengineering a commercial product
- first approach:  
split Linux' IP-Sec

# First Approach

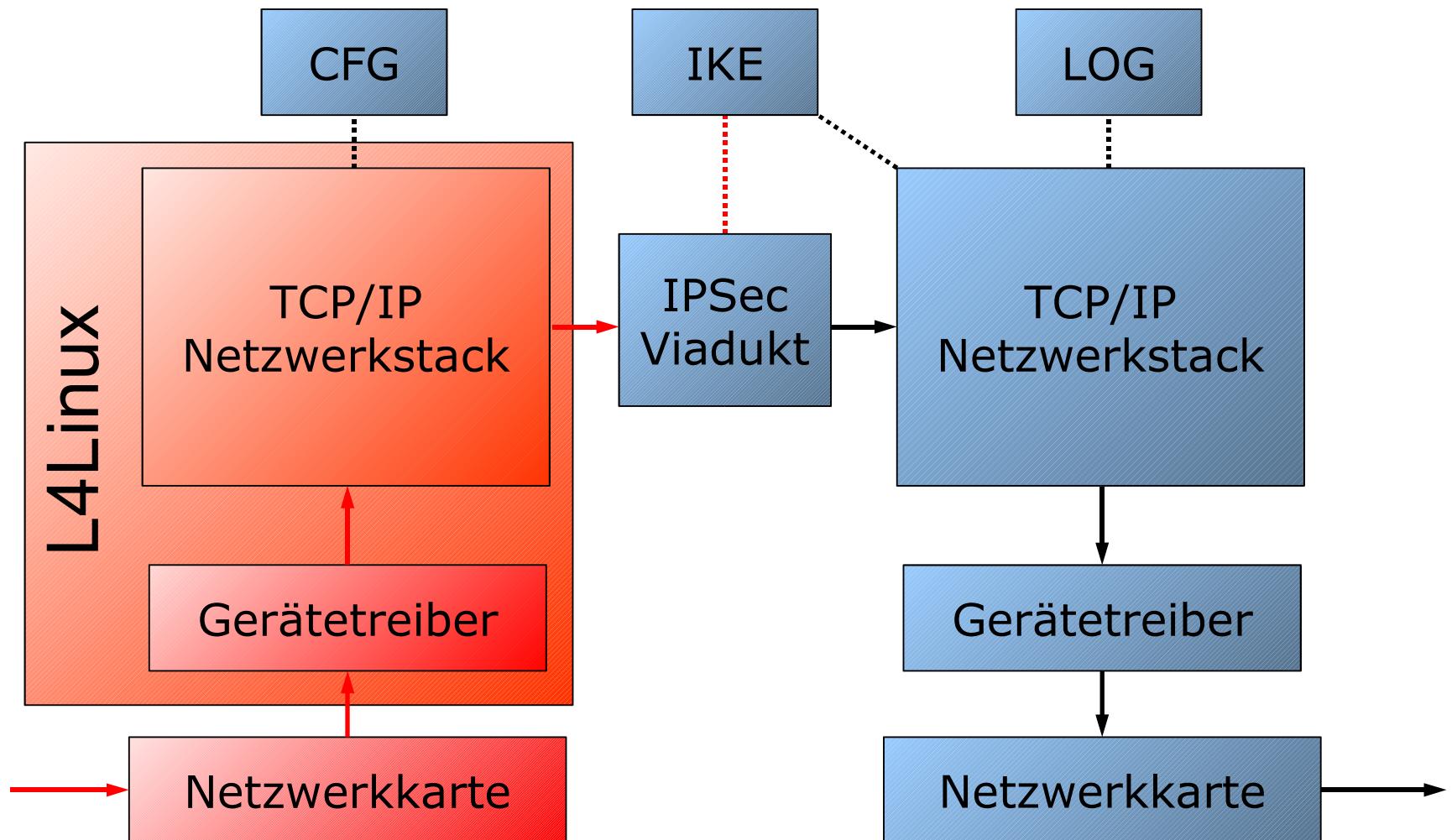
- Decompose network stack software on the basis of IP packet flow

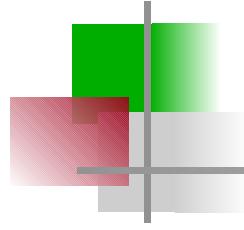


# Given up

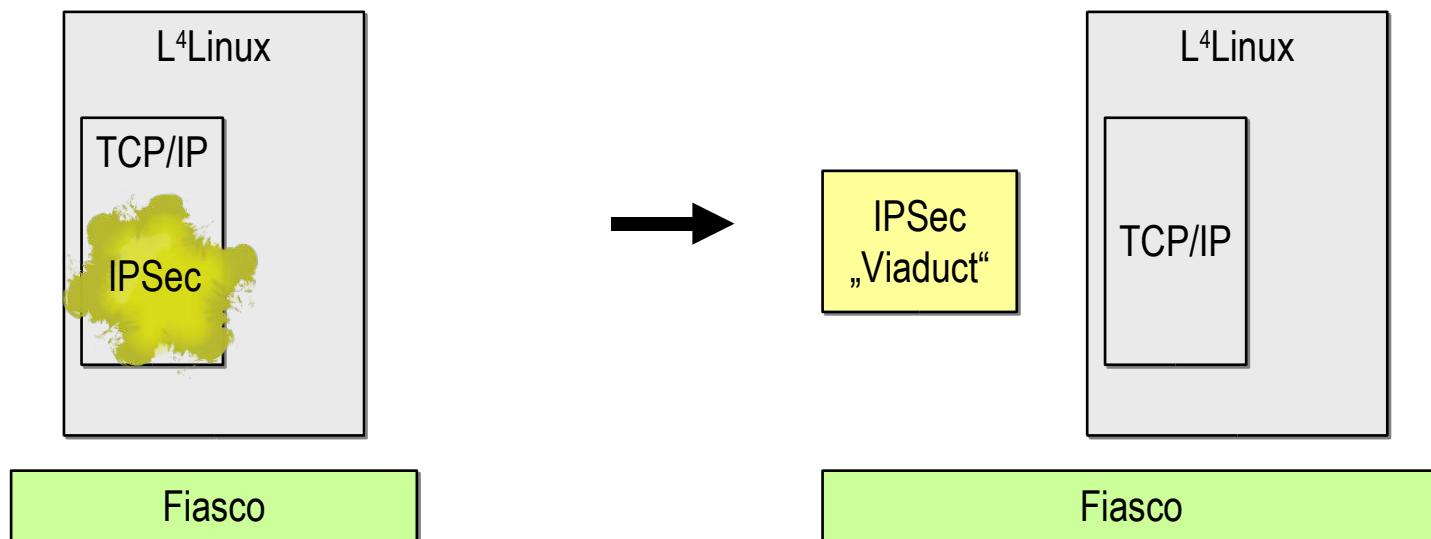


# Instead: Rewrite and Trusted Wrappers



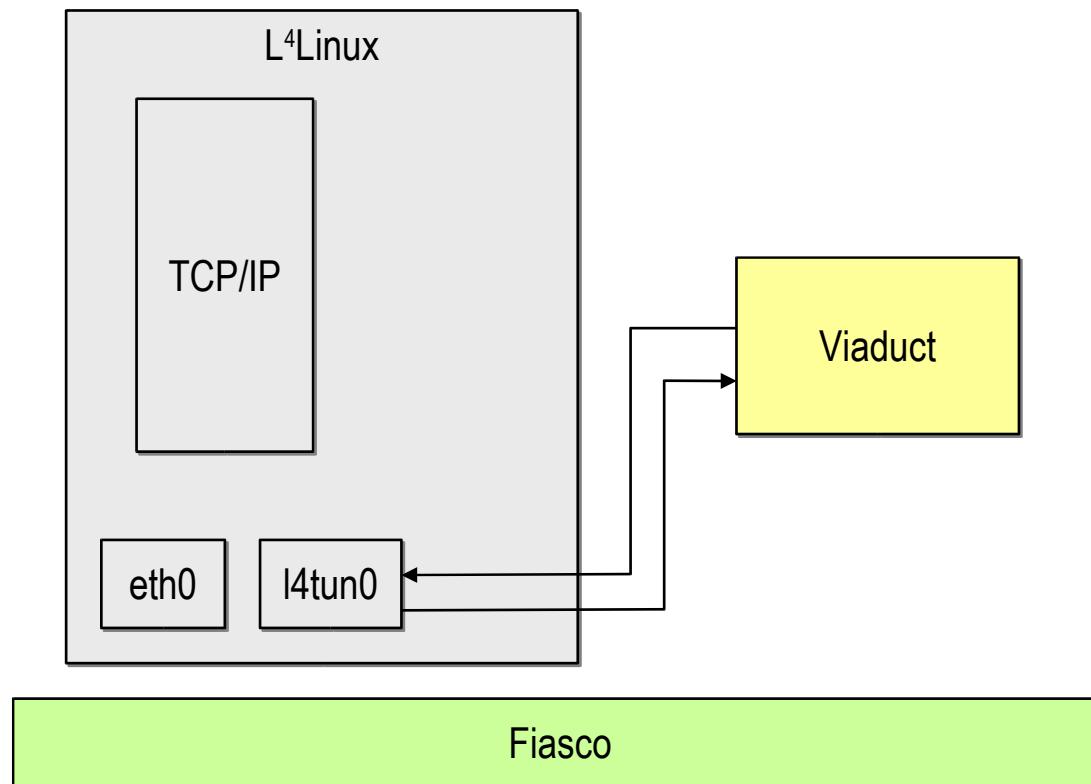


# One Step Back



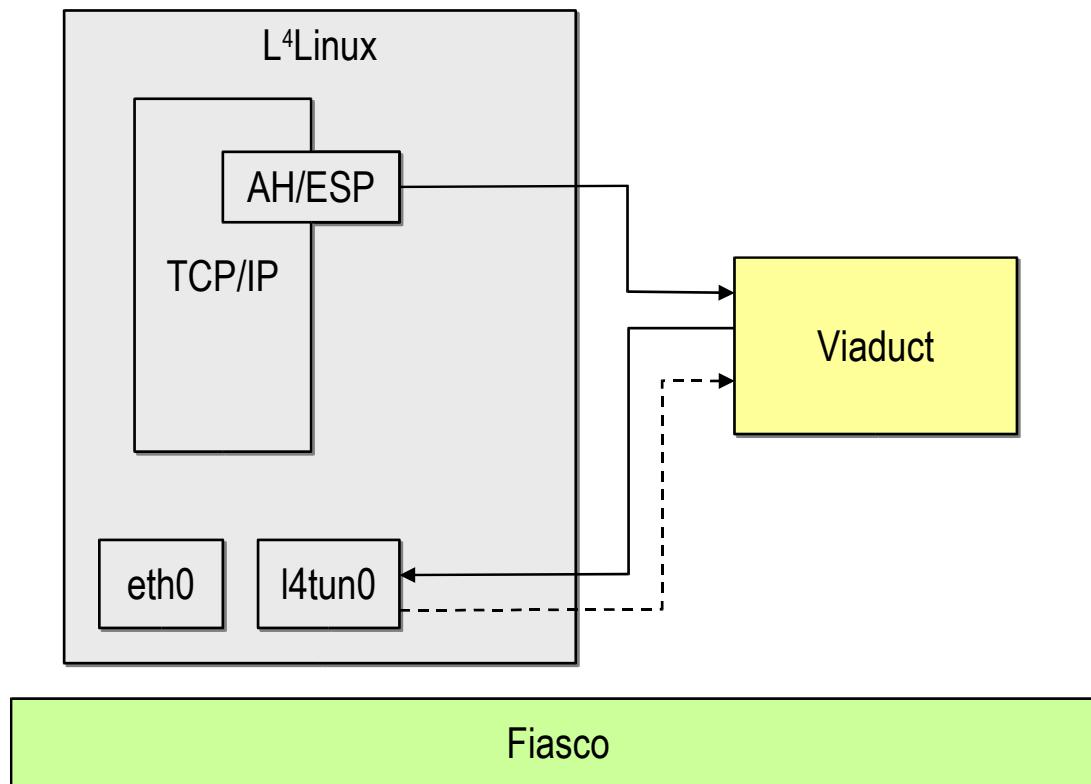
# Technical Details (Viaduct)

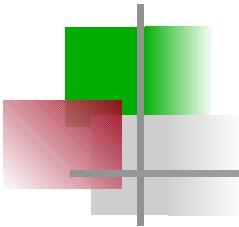
- IP packets must be passed to the Viaduct
  - L4 IPC as Virtual network driver in L<sup>4</sup>Linux



# Technical Details (Viaduct)

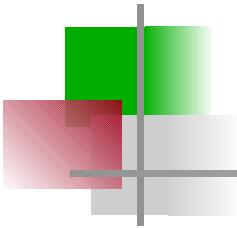
- IPSec can only handle unfragmented packets
  - Use L<sup>4</sup>Linux for complex reassembly





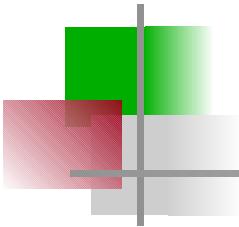
# Related Work: EROS, Keykos and Nizza

- similar objectives
- moving target



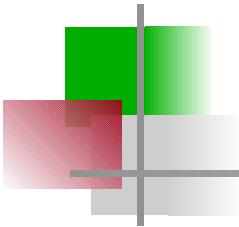
# Related Work: Microsoft NGSCB and Nizza

- similar objectives
- moving target



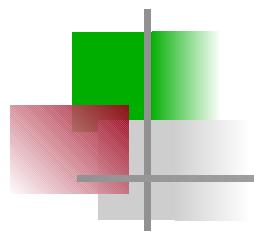
## Related Work: XOM and Nizza

- XOM: take OS off trusted path
- Implementing an Untrusted Operating System on Trusted Hardware  
David Lie Chandramohan A. Thekkath Mark Horowitz  
SOSP 2003



## Related Work: Terra and Nizza

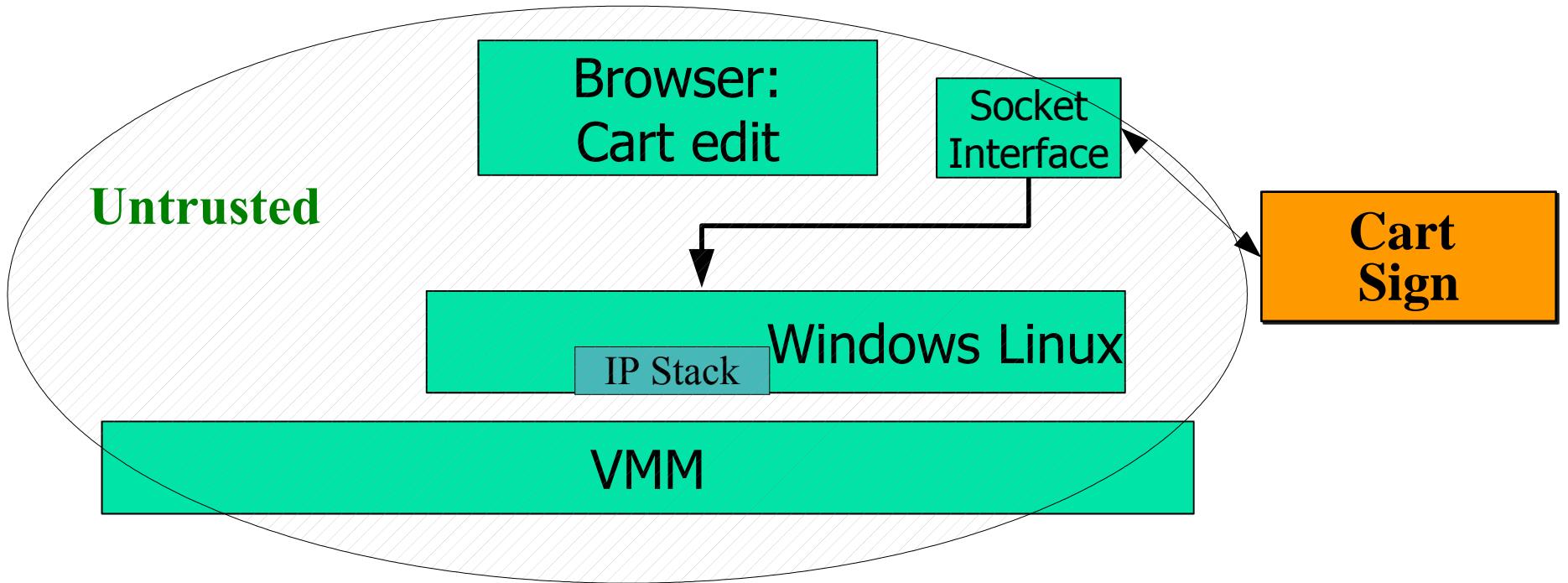
- Terra: VMM as trusted platform
- many more projects down that line
- A Virtual Machine-Based Platform for Trusted Computing  
Tal Garfinkel, Ben Pfaff, Jim Chow, Mendel Rosenblum and Dan Boneh  
SOSP 2003



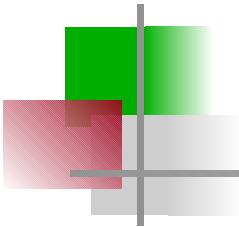
# Nizza vs. VMM approaches

- advantages VMM:
  - support+reuse: unmodified legacy OS
  - anything else ?
- advantages Nizza
  - smaller TCP:  
no network device emulation and drivers
  - fine grained sharing
  - efficiency (optimized message passing)
- VMM untrusted on Nizza ?

# Challenge: Untrusted VMM

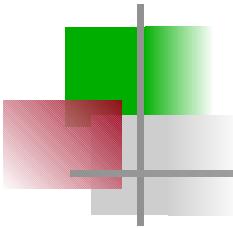


Minimal Trusted Platform



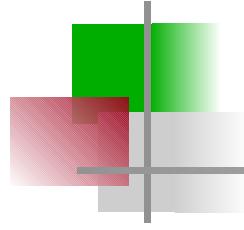
# Related work “Useful” for Nizza

- Secure Storage on Untrusted Servers Secure Untrusted Data Repository (SUNDR)  
Jinyuan Li, Maxwell Krohn, David Mazire s,  
and Dennis Shasha ,  
New York University
- Privtrans: Automatically Partitioning Programs for Privilege Separation  
David Brumley and Dawn Song,  
Carnegie Mellon University



# Technical Risks

- performance
  - copying overhead
  - context switching time (hardware)
  - increased memory  
(duplication of page tables)



# Context switches

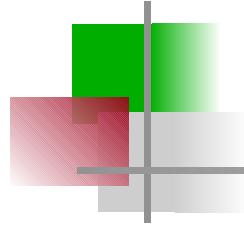
- Register IPC between two address spaces  
(1 x send, 1 x receive; kernel entry with sysenter):

Pentium-III: 600 cycles

Opteron: 700 cycles

Prescott: 2200 cycles

??



# Conclusion

technologies are in place  
to build much better (securer) systems

need proper integration → Nizza

