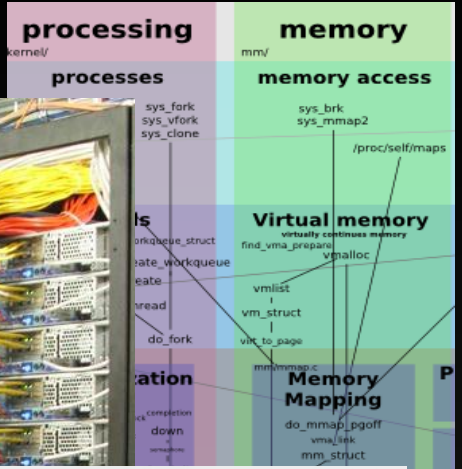# XenTT: Deterministic Systems Analysis in Xen

Anton Burtsev, David Johnson, Chung Hwan Kim, Mike Hibler, Eric Eide, John Regehr
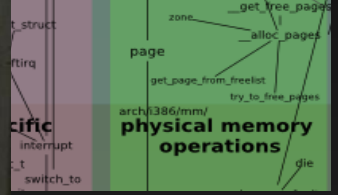
University of Utah

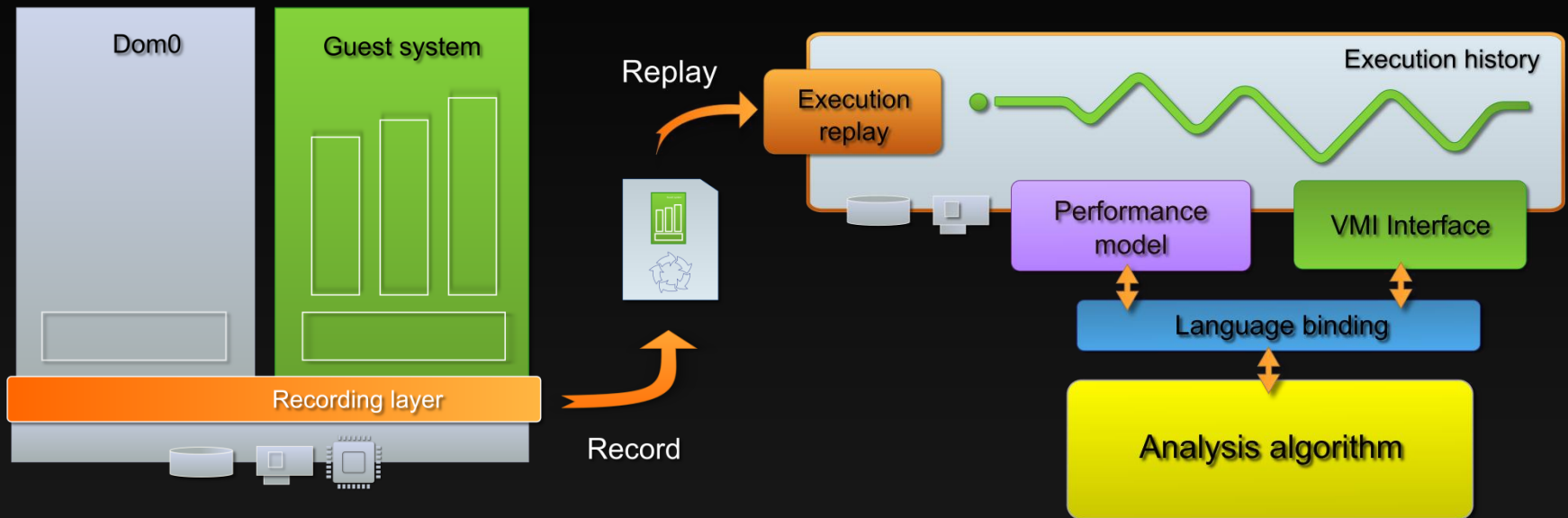August 27-28, 2012
San Diego, CA, USA

- Record execution history of a guest VM
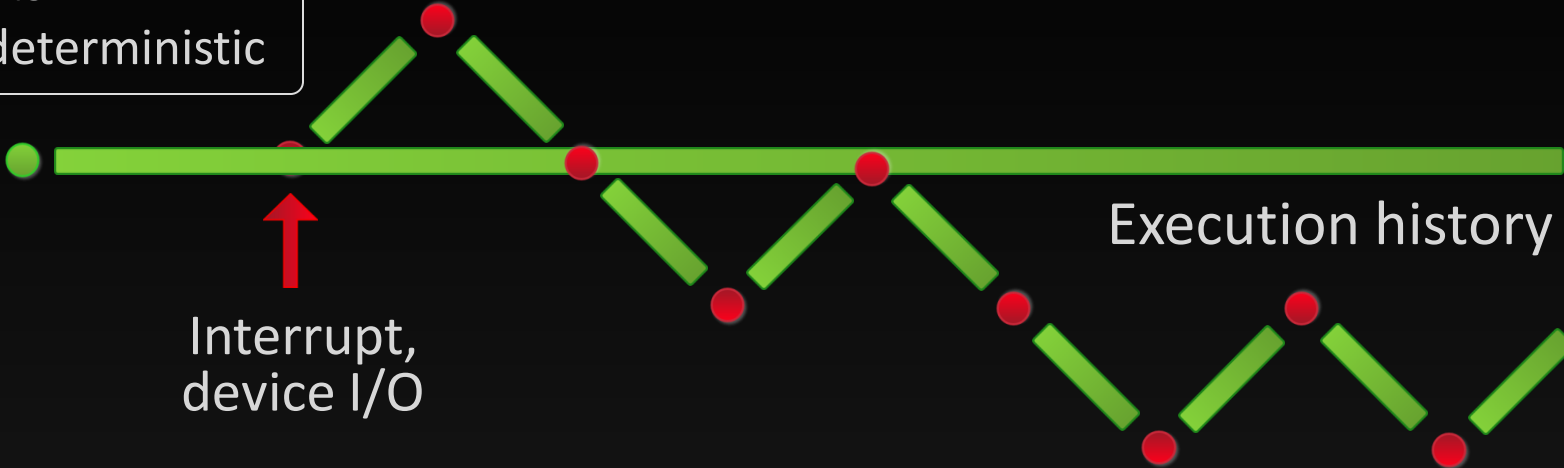- Recreate it in an instruction-accurate way

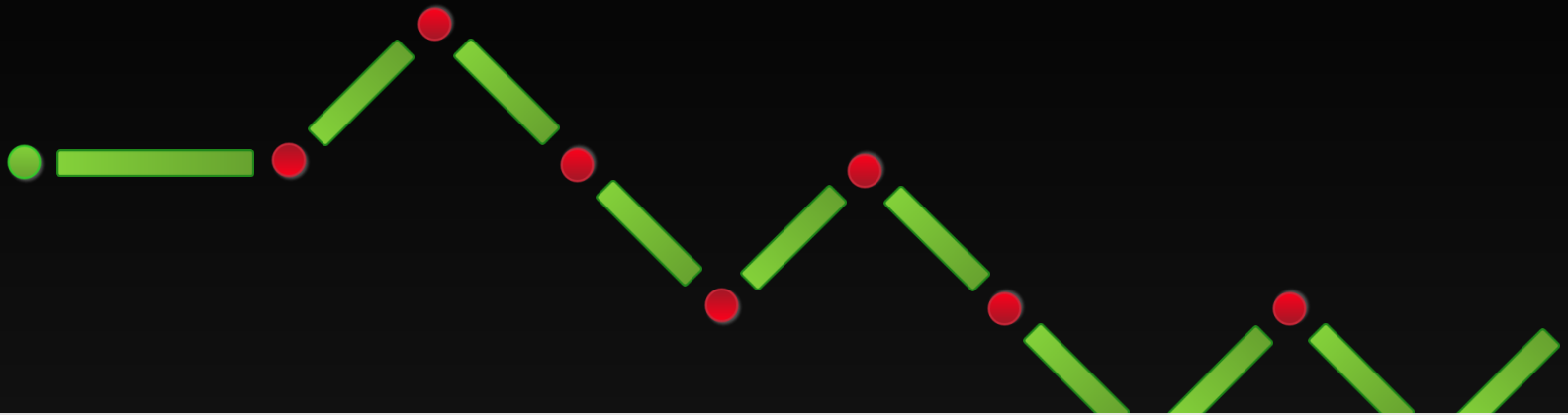- Execution replay is a right way to analyze systems
- We need a practical tool!

# Deterministic replay

# Determinism

CPU is deterministic

Execution history

Interrupt, device I/O

# Recording



- Determinism of the execution environment
- Instruction-accurate position of events

Event log

# Instruction-accurate position of events

```
label:    ...
          mov
          shr
          mov
➡    rep  movsl
          test
          jne  label
          ...
```

- Number of instructions since boot
  - Intel has a hardware counter
  - It's not accurate

- Hardware instruction counter
  - Preempt execution of a system at the same instruction
  - Hardware instruction counter + single-stepping

# Determinism in Xen

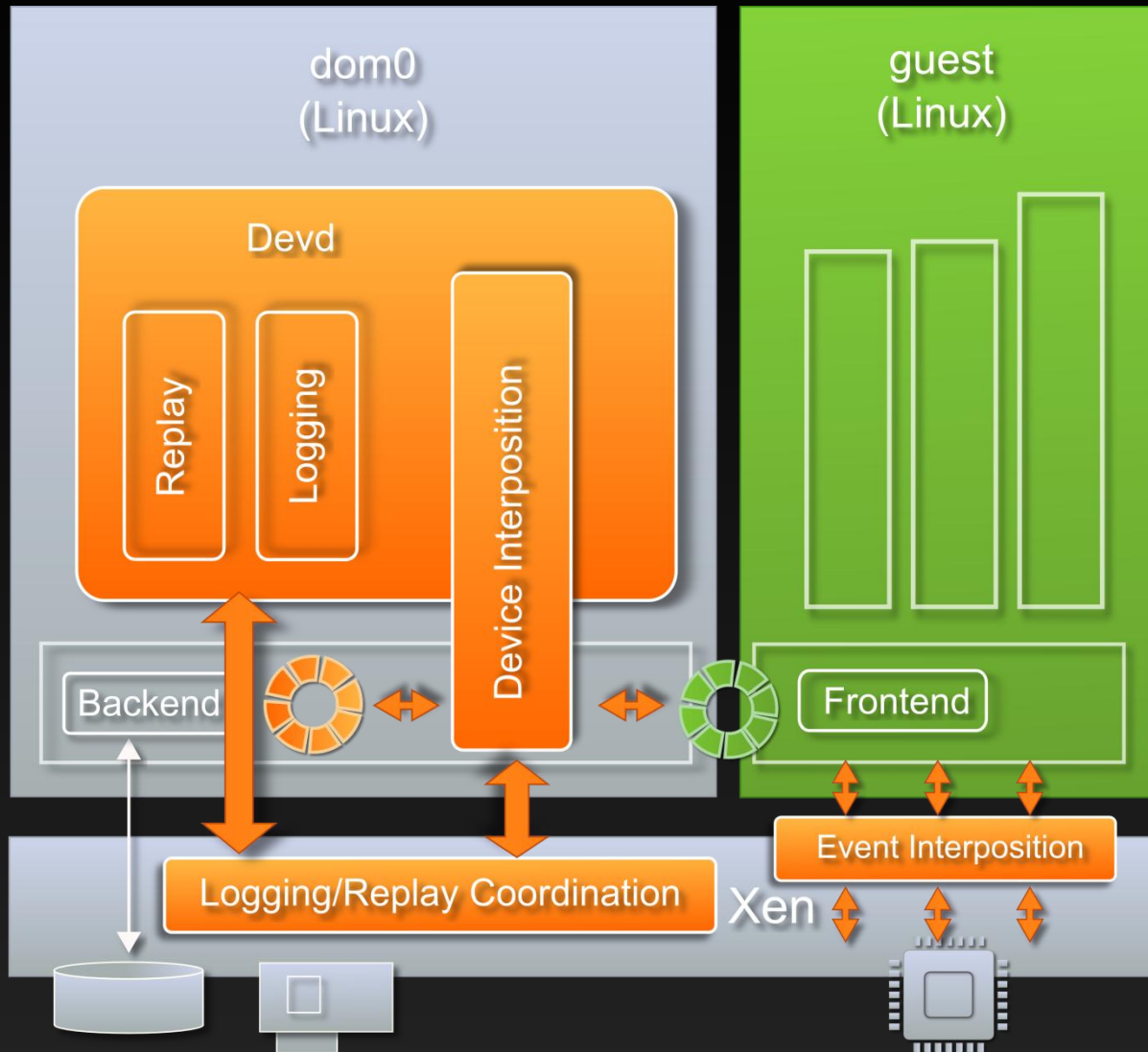# Nondeterministic events

Simple Model

System = memory pages + registers

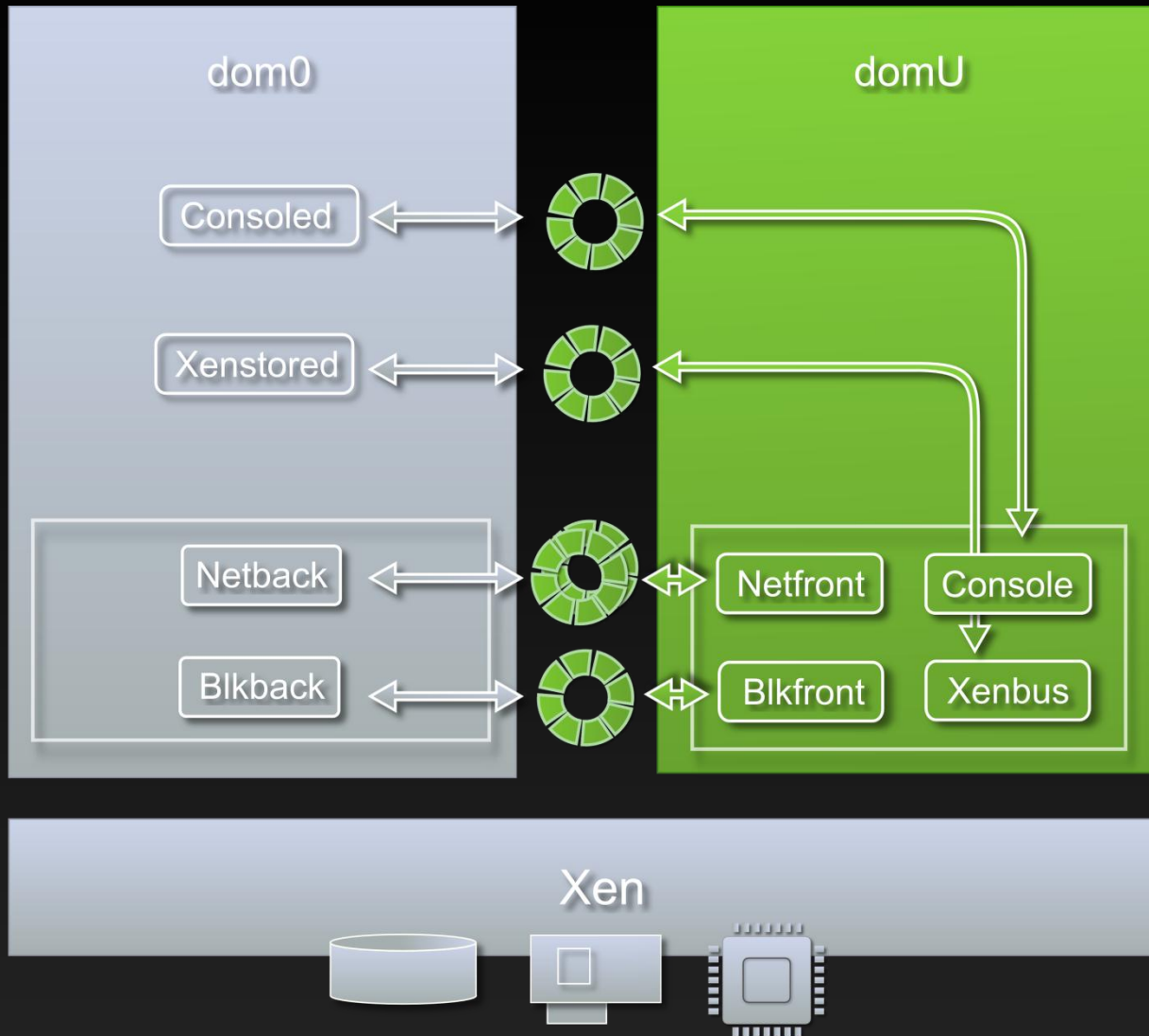Events = memory updates
(time, device I/O, system calls)

Control flow updates = registers + stack
(interrupts, events)
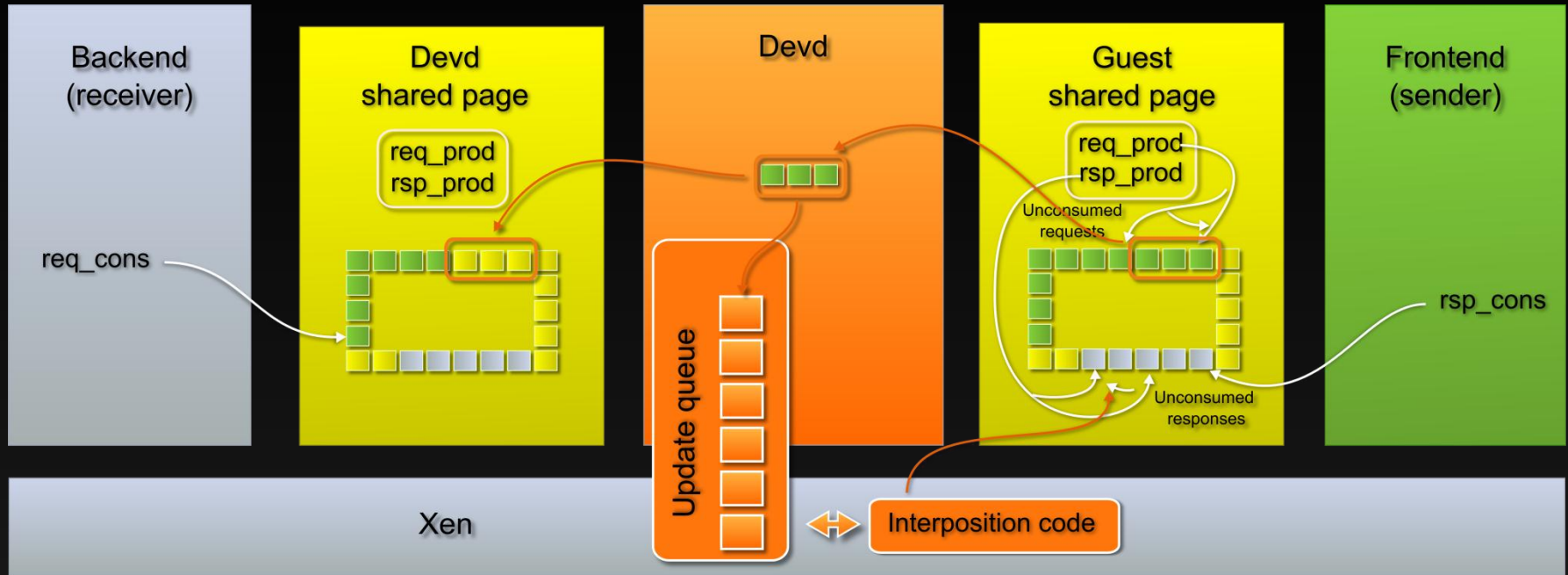
# Some examples

- Instruction emulation (e.g. cpuid, rdtsc, in/out)
  - Return the values of the original run
- Hypercalls
  - Re-execute to ensure determinism of the hypervisor
- Time
  - Shared info page + rdtsc
- Exceptions
  - Deterministic, re-execute
- Interrupts
  - Force re-execution of the interrupt frame (bounce frame) code in entry.S
- Shared info updates
  - Replay original values
- Memory
  - Shadow page tables

11

# Xen devices
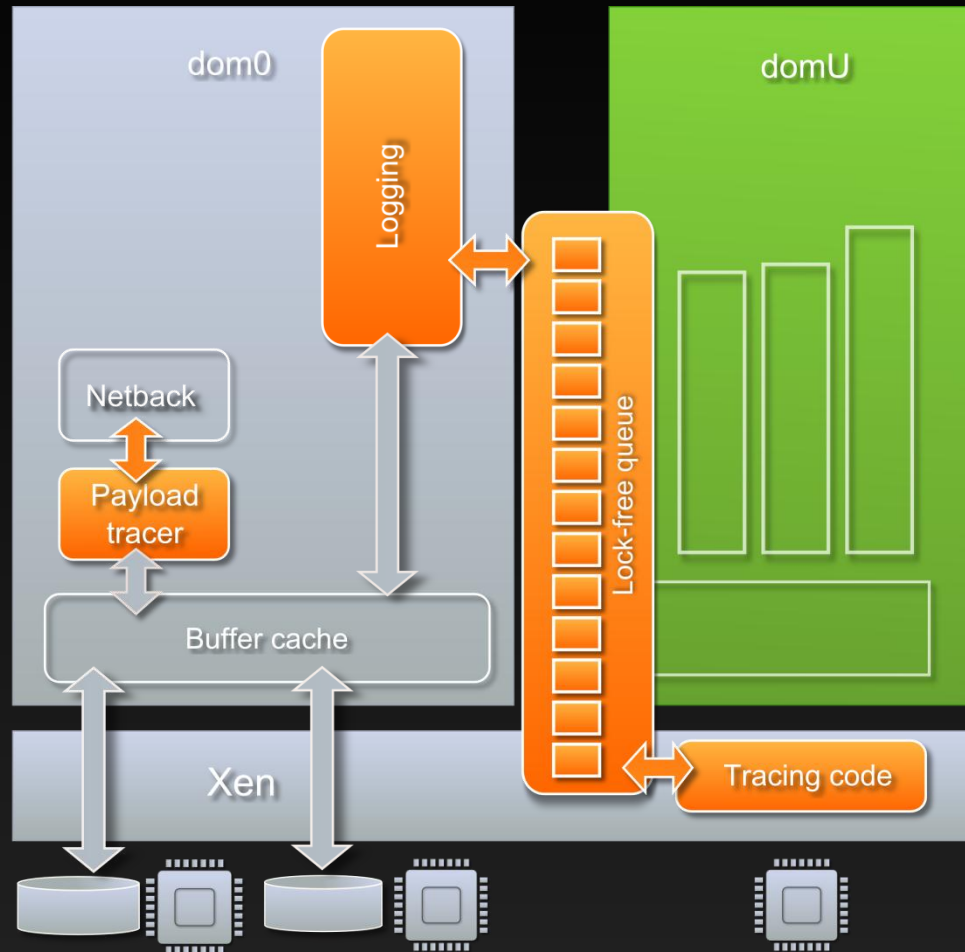
# Device interposition



- Devd ensures determinism of updates to the guest's shared ring buffers

# Replay touches many parts of Xen

- Device discovery
  - Devd implements a concept of a device bus
  - Discovers new devices in Xenstore
  - Binds new devices with drivers

- Xenstore transactions
  - During replay, transactions from replayed guest can't fail
  - They will not be re-executed

- Out-of-order device responses
  - Disk and network responses can arrive out-of-order

- Disk logging
  - Disk payload is deterministic
  - LVM snapshots

- Network logging
  - In-kernel logging of the network payload

# Low-overhead logging
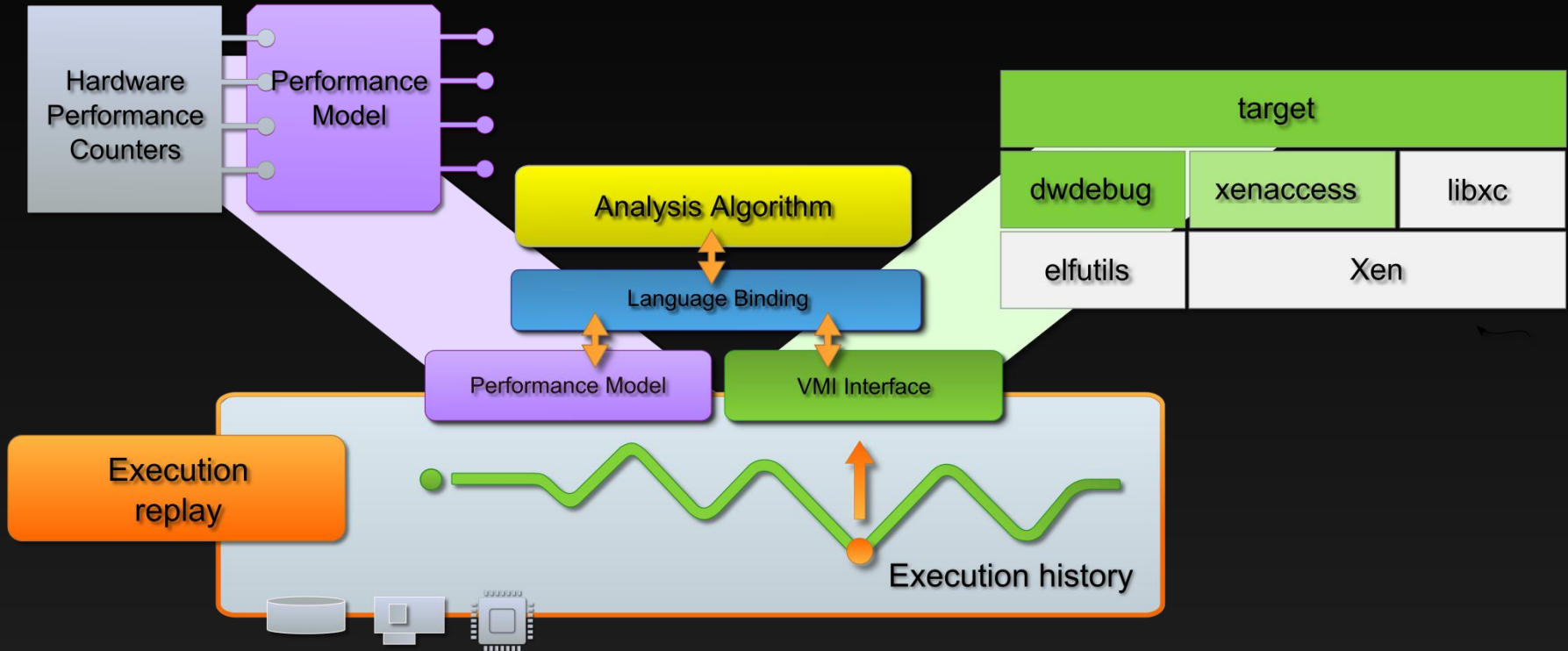
# Are we sure that executions identical?

- Intel branch store trace facility
- Record all taken branches in a memory buffer
  - Compare original and replay runs

# Analysis Engine and
# Virtual Machine Introspection

# Analysis framework

# Virtual Machine Introspection (VMI)

# Performance model



Hardware CPU

Performance Model

○──● Hardware performance counters

○──● Virtual performance counters

- Account for effects of replay

- Translate performance between original and replay runs

- Re-execution approach to performance

$$\text{Virt cntr} = \text{Virt cntr}_{start} + \Delta\,(\text{Real cntr})$$

# Analysis Examples

# NFS request processing path



- How much time requests spend in each subsystem?

- Request tracking
  - Address of the kernel data structure is a unique identifier
  - Join identifiers when requests move between subsystems

22

# Execution context tracking

- Execution context
  - Context switches
    - `schedule.switch_tasks`
  - User/kernel
    - System call transitions
    - `system_call`
  - Interrupts and exceptions
    - `do_IRQ`
    - `do_pagefault`
    - `do_*` (`divide_error, debug, nmi, int3...`)

- Make analysis context aware
  - Filter probes by context, e.g.
  - All pagefaults from the process "foo"

# Control Flow Integrity (CFI)

Disassemble
the entry point

Static
disassemble
fails?

Yes

Single-step
until next entry point

Run until
next probe

- Simple CFI model:
  - Returns should match calls
  - Dynamic calls are "sane", e.g. within kernel address space
- Detect ROP attacks, stack smashing, etc.

# Execution trace

```
sys_sendfile
   do_sendfile
      fget_light
      rw_verify_area
      fget_light
      rw_verify_area
      shmem_file_sendfile
         do_shmem_file_read
            shmem_getpage
               find_lock_page
                  radix_tree_lookup
               shmem_recalc_inode
               shmem_swp_alloc
                  shmem_swp_entry
                     kmap_atomic
                        __kmap_atomic
                        page_address
               kunmap_atomic
               find_get_page
                  radix_tree_lookup
            file_send_actor
            sock_sendpage
               UNKNOWN FUNCTION
               (addr:0x00000000)
```

- CFI records a trace of function calls
  - sys_sendfile is the last system call before control flow jumps to 0x0

# Intrusion backtracking



write /etc/passwd

- Track accesses to "/etc/passwd"
- Probe sys_open
  - Filter by file name
  - Find process ID, branch counter

# Intrusion backtracking:  pass 1

raise current
pid privilege

write /etc/passwd

- Process or it's parent escalated privileges
- Watch write accesses  to &task->uid
  - Filter by parents of the offending process

# Intrusion backtracking:  pass 1

raise current
pid privilege

write /etc/passwd

last system call

- Find the syscall inside which privileges are escalated
- Probe sys_* - all system call entry and exit points
  - Filter by the offending process ID

28

# Intrusion backtracking:  pass 1



raise current pid privilege

kernel jumps to 0x0

write /etc/passwd

last system call

- Privilege escalation is a CFI violation
- Start CFI analysis from the last system call
- Find %EIP at which CFI is violated, and location of the shell code (0x0)

# Intrusion backtracking:  pass 1



- Find at which point address 0x0000000 gets mapped
- Probe do_page_fault

# More mechanisms

- Execution traces
  - BTS trace of all taken branches
  - Instruction traces

- Memory (variable) access traces
  - Intersect with the execution trace
  - See where variables get accessed

# How much overhead?

- 32-bit x86 PV-guests
- xen-unstable near v3.0.4
  - We rely on a working shadow page tables
- 1-CPU time-traveling guests
  - No SMP replay
  - Dom0 and Xen are SMP of course

- Test machine
  - 4 cores
  - 1Gbps network
  - 130 MB/s disks

# CPU

# Network throughput

# Network delay



Ping

Legend: Xen (green), XenTT (orange)

# Disk I/O

|  | **Raw** | **Compressed (gzip)** |
|---|---|---|
| **Linux boot** | | |
| Event log | 4.3 GB | 0.8 GB |
| **Idle overnight (14 hours)** | | |
| Event log | 6.2 GB | 1.6 GB |
| Growth rate | 440 MB/hour | 114 MB/hour (2.7GB/day) |
| **TCP receive (1.63 GB data stream)** | | |
| Event log | 1.76 GB | 342 MB |
| Payload log | 1.79 GB | Payload dependent |
| **Disk write (4 GB file)** | | |
| Event log | 1.8 GB | 350 MB |
| **Disk read (4 GB file)** | | |
| Event log | 0.59 GB | |

# Lessons learnt

# Scaling development

- Extending Xen with BTS support
  - Debug crashes in Xen, and dom0

- Execution comparison tools
  - BTS traces to understand what went wrong
  - Support for resolving symbols

- Run-time comparison tools
  - Compare guest's state between original and replay runs

- Trace from all parts of your system
  - Xen, dom0, domU

- Support performance tracing
  - Xentrace messages

# What we didn't predict

- I/O delay goes up
  - Not sure if Linux has adequate low-latency user-level processing support
  - Maybe need an in-kernel interposition component

- Branch counters are fragile
  - Our code works on several server CPUs
  - Fails on a laptop with the CPU from the same model/family line

# Conclusions

# Practical replay analysis is feasible

- Performance overheads are reasonable
  - Realistic systems
  - Realistic workloads
  - Minor setup costs (just install Xen)

- Analysis engine is an amazing tool
  - And we're growing it

- We need your help to port it upstream
  - Porting effort
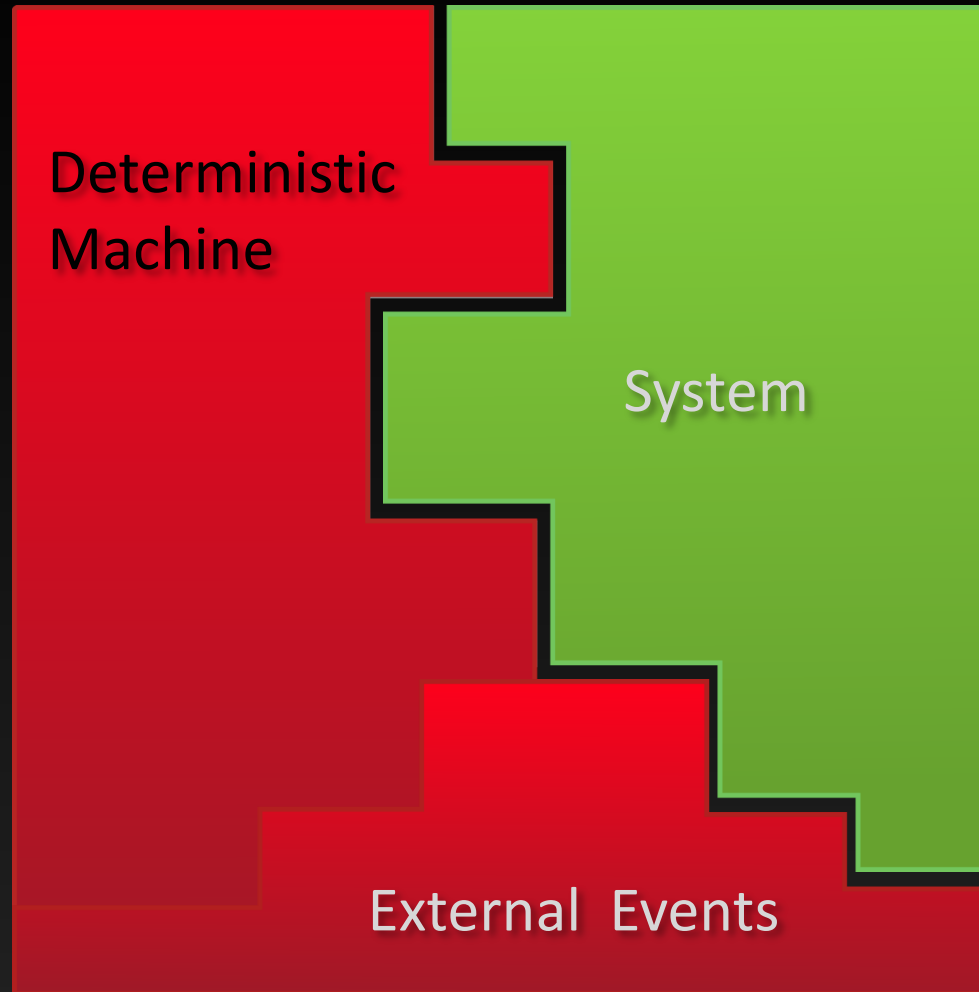  - Shadow memory for PV guests
  - Support for HVM guests

# Thank you.

All code is GPLv2 and will be available soon
(available on-request now).

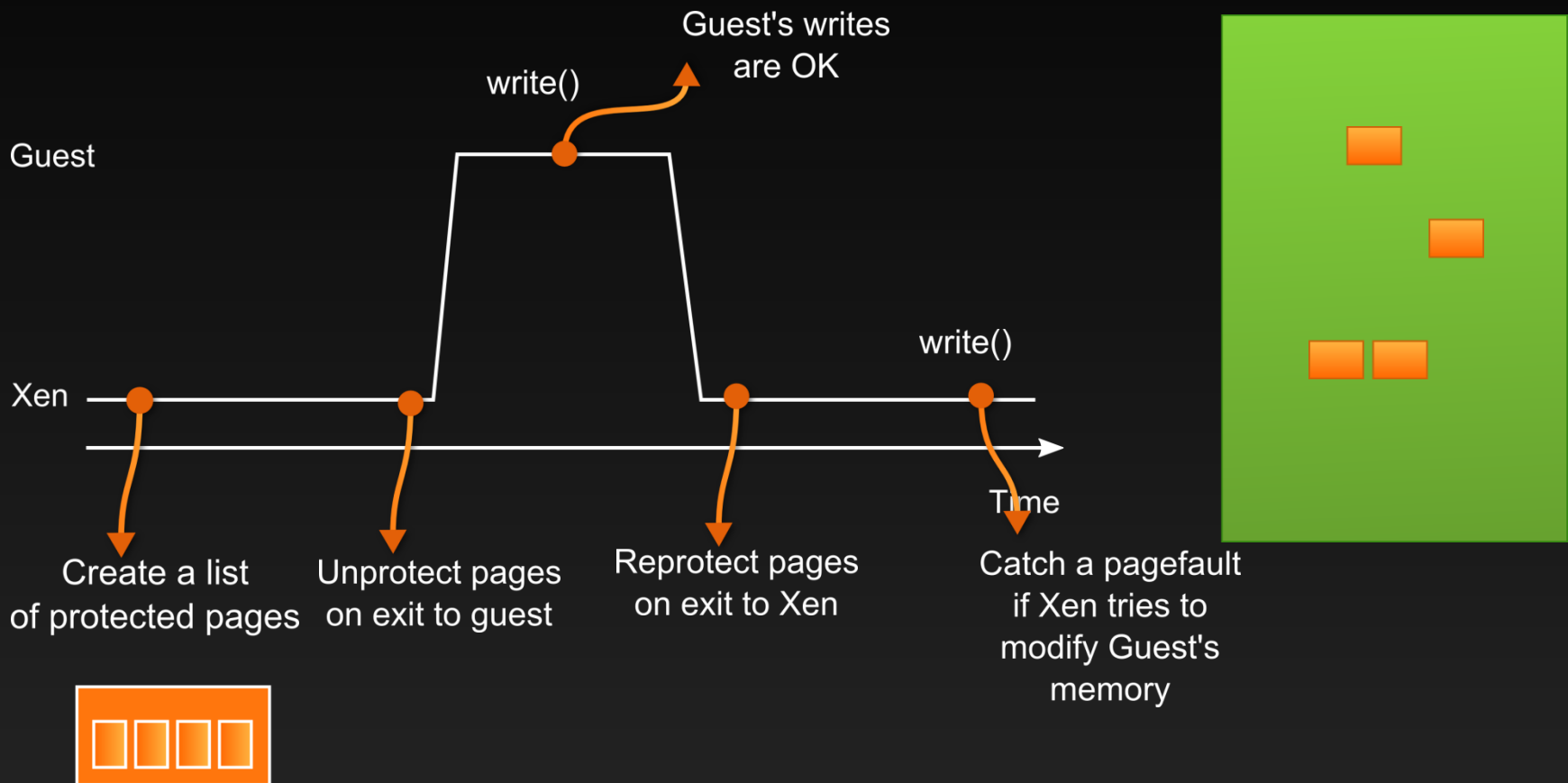Questions: aburtsev@flux.utah.edu

# Backup slides

# Execution environment
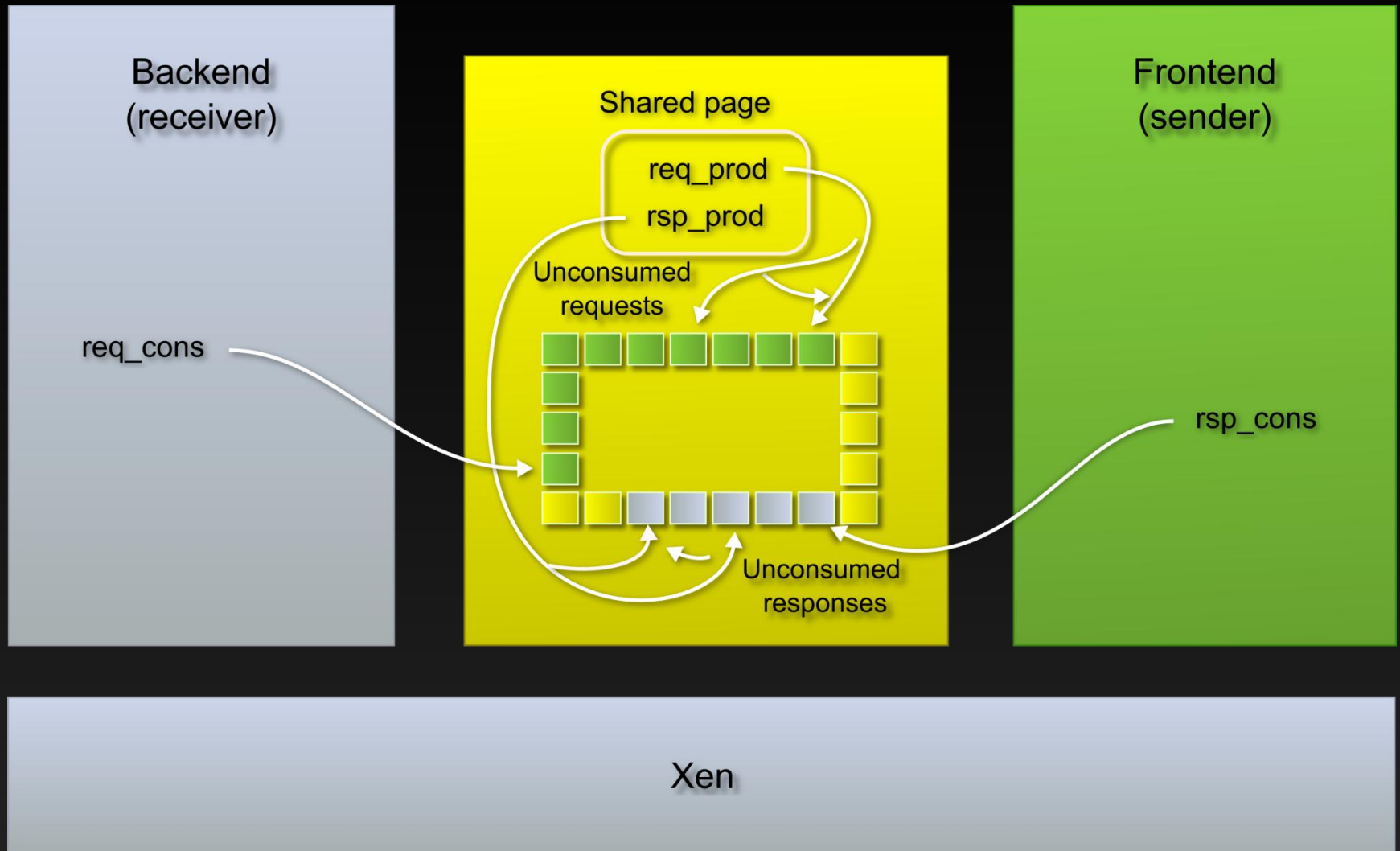
- CPU
- Virtual hardware
  - Memory
  - Disk
- Software

Deterministic Machine

System

External Events

- External I/O

# How do we find nondeterministic events?

# Device interface

# Reading a local variable

```
bsym_skb = target_lookup_sym(probe->target,
                             "netif_poll.skb",
                             ".", NULL, flags);

lval_skb = bsymbol_load(bsym_skb, flags);

skb_addr = *(unsigned long*) lval_skb->buf;
```

# Controlled re-execution or replay

- Types of non-deterministic events
  - Synchronous
    - Hypercalls
  - Asynchronous
    - Interrupts
  - Best effort
    - Time updates
- Branch counters
  - The biggest problem of this solution