COMP9242
Advanced Operating Systems

S2/2013 Week 8:
Virtualization

# Copyright Notice

**These slides are distributed under the Creative Commons Attribution 3.0 License**

- You are free:
  - to share—to copy, distribute and transmit the work
  - to remix—to adapt the work

- under the following conditions:
  - **Attribution:** You must attribute the work (but not in any way that suggests that the author endorses you or your use of the work) as follows:
    - "Courtesy of Gernot Heiser, [Institution]", where [Institution] is one of "UNSW" or "NICTA"

The complete license text can be found at
http://creativecommons.org/licenses/by/3.0/legalcode

# Virtual Machine (VM)

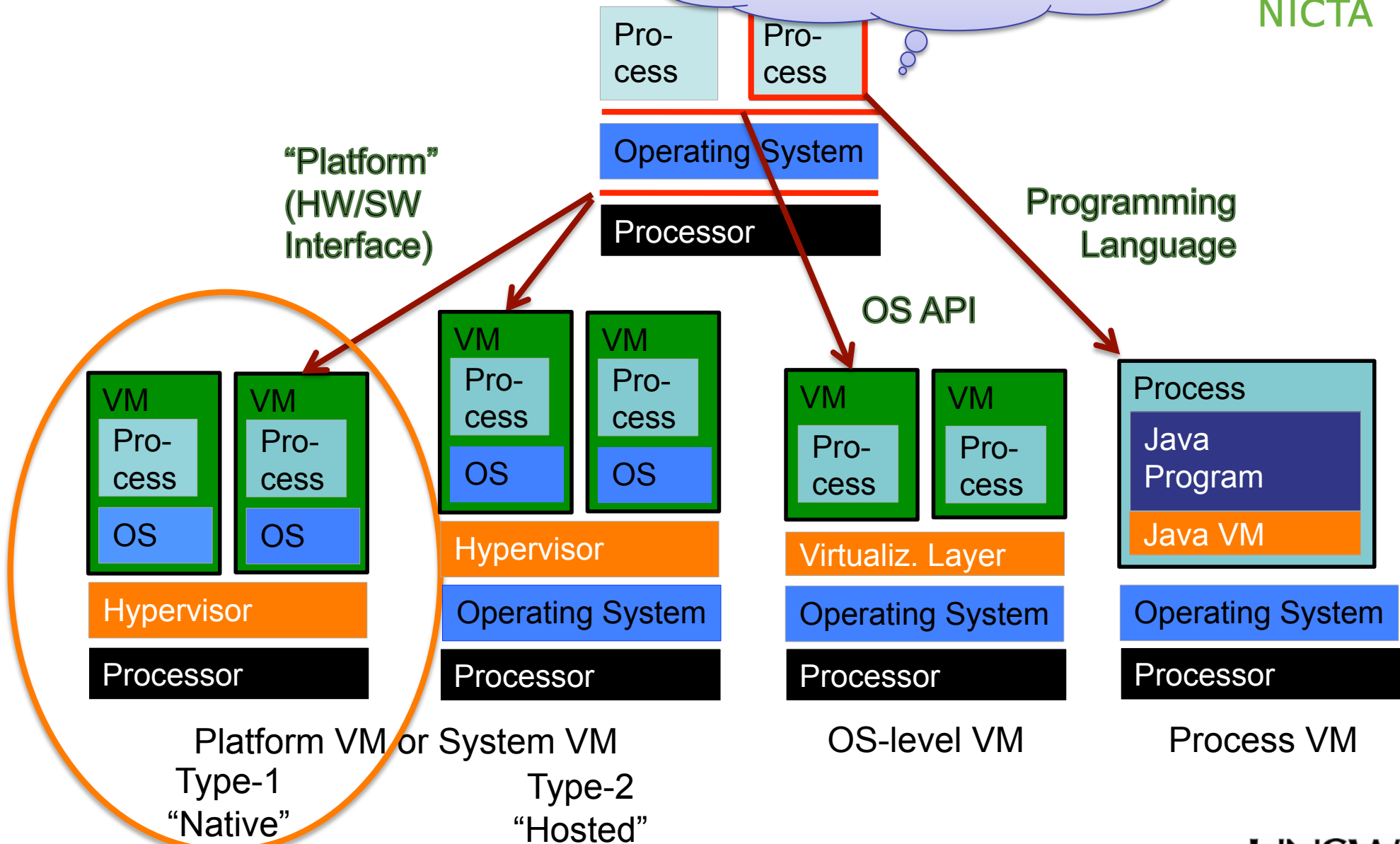*"A VM is an efficient, isolated duplicate of a real machine"*

- Duplicate: VM should behave identically to the real machine
    - Programs cannot distinguish between real or virtual hardware
    - Except for:
        - Fewer resources (and potentially different between executions)
        - Some timing differences (when dealing with devices)
- Isolated: Several VMs execute without interfering with each other
- Efficient: VM should execute at speed close to that of real hardware
    - Requires that most instruction are executed directly by real hardware

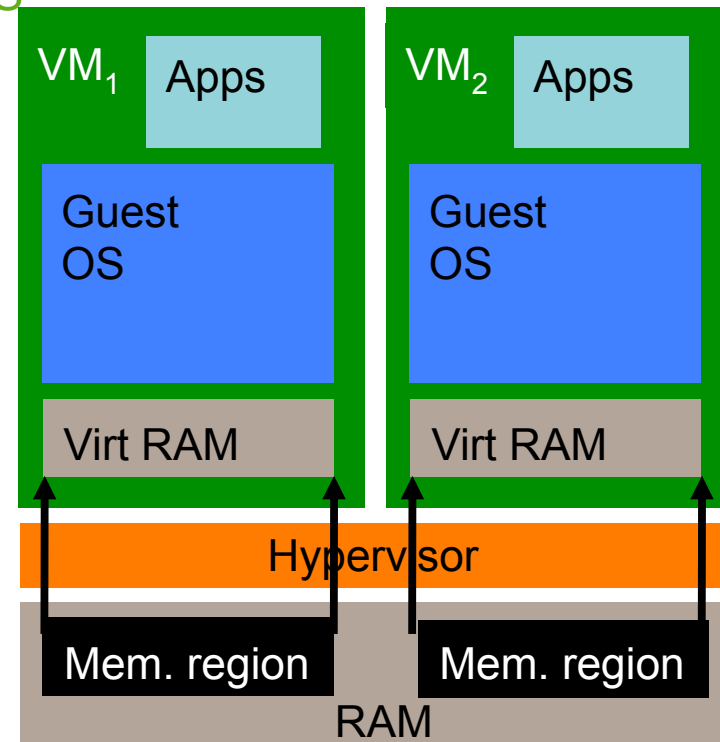*Hypervisor* aka *virtual-machine monitor*: Software implementing the VM

NICTA

UNSW
THE UNIVERSITY OF NEW SOUTH WALES

# Types of Virtualization

Plus anything else you want to sound cool!

NICTA

Process

Process

"Platform" (HW/SW Interface)

Operating System

Processor

Programming Language

OS API

| VM | VM |
|---|---|
| Process | Process |
| OS | OS |

Hypervisor

Processor

Platform VM or System VM
Type-1 "Native"

| VM | VM |
|---|---|
| Process | Process |
| OS | OS |

Hypervisor

Operating System

Processor

Type-2 "Hosted"

| VM | VM |
|---|---|
| Process | Process |

Virtualiz. Layer

Operating System

Processor

OS-level VM

Process

Java Program

Java VM

Operating System

Processor

Process VM

UNSW
THE UNIVERSITY OF NEW SOUTH WALES

# Why Virtual Machines?

- Historically used for easier sharing of expensive mainframes
  - Run several (even different) OSes on same machine
    - called *guest operating system*
  - Each on a subset of physical resources
  - Can run single-user single-tasked OS in time-sharing mode
    - legacy support
- Gone out of fashion in 80's
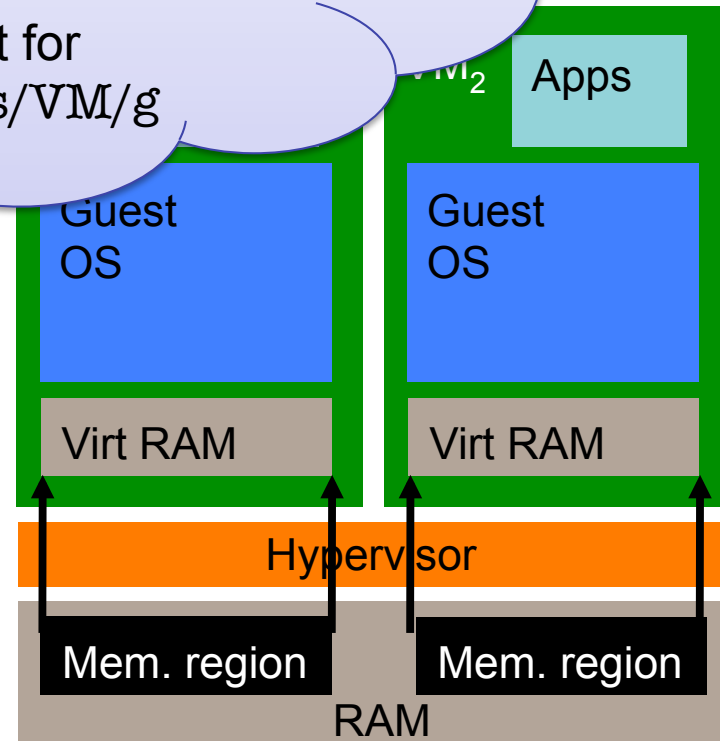  - Time-sharing OSes common-place
  - Hardware too cheap to worry...

# Why Virtual Machines?

- Renaissance in recent years for improved isolation
- Server/desktop virtual machines
  - Improved QoS and
  - Uniform view of
  - Complete enc
    - replication
    - migration
    - checkpointing
    - debugging
  - Different concurrent OSes
    - eg Linux + Windows
  - Total mediation
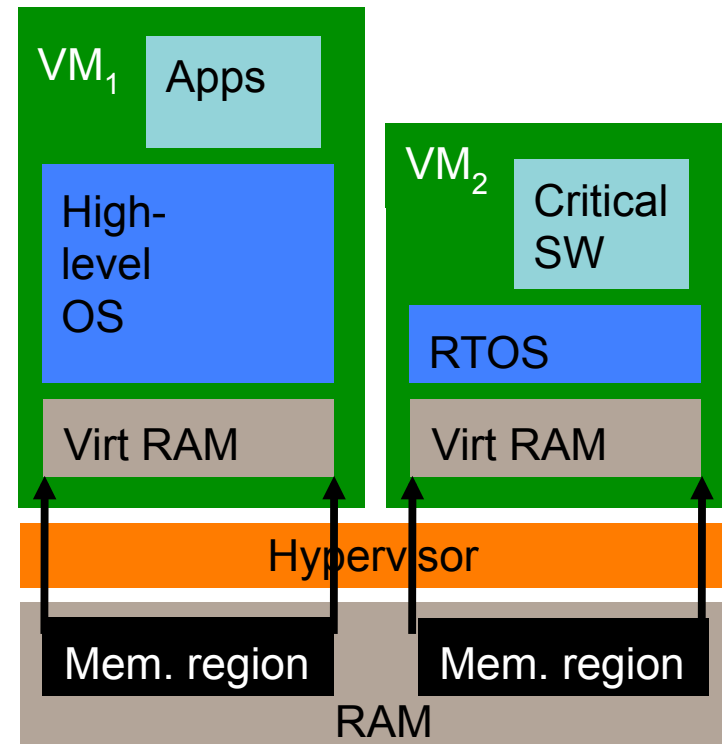- Would be mostly unnecessary
  - … if OSes were doing their job!

Gernot prediction of 2004:
2014 OS textbooks will be
identical to 2004 version
except for
s/process/VM/g

Apps

Guest OS

Guest OS

Virt RAM

Virt RAM

Hypervisor

Mem. region

Mem. region

RAM

NICTA

UNSW
THE UNIVERSITY OF NEW SOUTH WALES
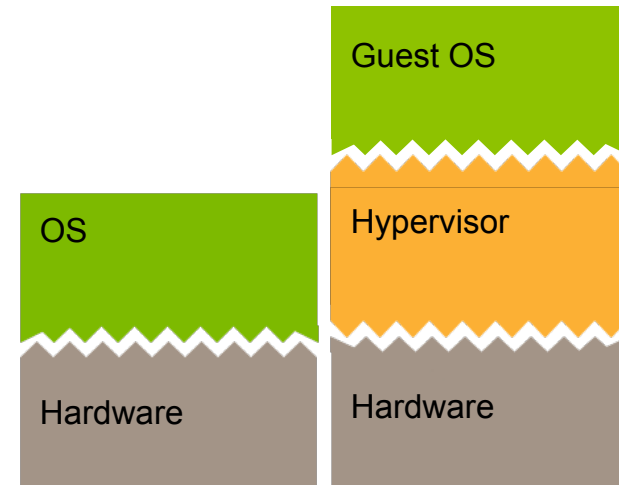
# Why Virtual Machines?

- Embedded systems: integration of heterogenous environments
    - RTOS for critical real-time functionality
    - Standard OS for GUIs, networking etc
- Alternative to physical separation
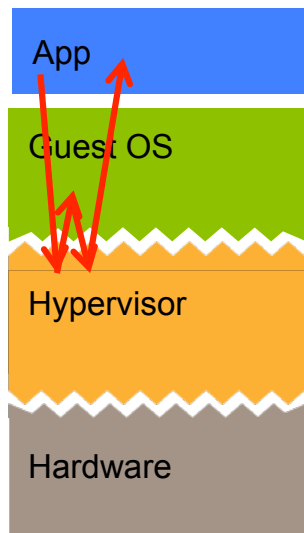    - low-overhead communication
    - cost reduction

# Hypervisor

- Program that runs on real hardware to implement the virtual machine
- Controls resources
  - Partitions hardware
  - Schedules guests
    - "*world switch*"
  - Mediates access to shared resources
    - e.g. console
- Implications
  - Hypervisor executes in *privileged* mode
  - Guest software executes in *unprivileged* mode
  - Privileged instructions in guest cause a trap into hypervisor
  - Hypervisor interprets/emulates them
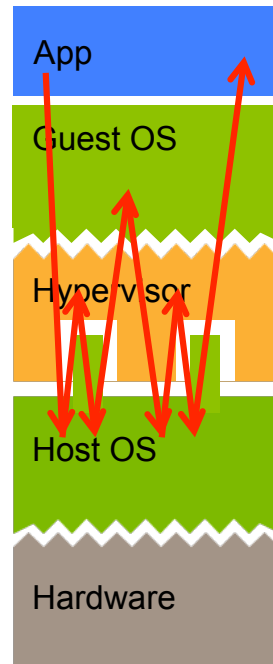  - Can have extra instructions for *hypercalls*

# Native vs. Hosted VMM

**Native/Classic/Bare-metal/Type-I**

| App |
| Guest OS |
| Hypervisor |
| Hardware |

**Hosted/Type-II**

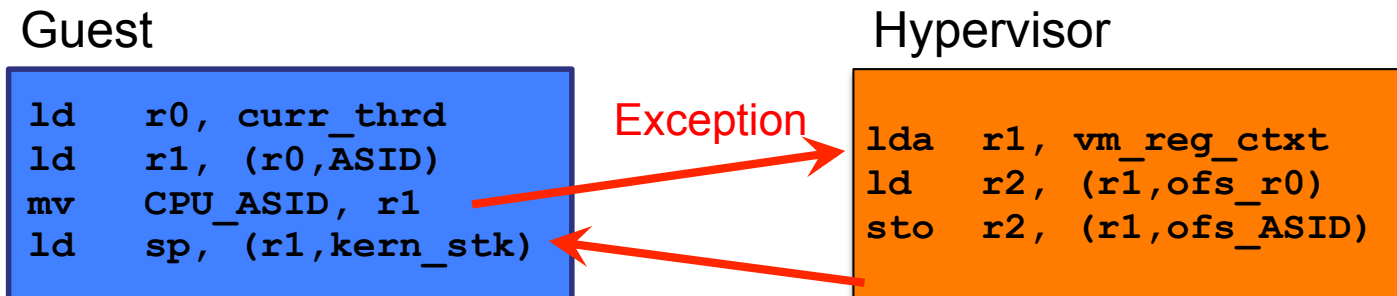| App |
| Guest OS |
| Hypervisor |
| Host OS |
| Hardware |

- Hosted VMM beside native apps
  - Sandbox untrusted apps
  - Convenient for running alternative OS on desktop
  - leverage host drivers

- Less efficient
  - Double node switches
  - Double context switches
  - Host not optimised for exception forwarding

# Virtualization Mechanics: Instruction Emulation

- Traditional *trap-and-emulate* (T&E) approach:
  - guest attempts to access physical resource
  - hardware raises exception (trap), invoking HV's exception handler
  - hypervisor emulates result, based on access to virtual resource
- Most instructions do not trap
  - prerequisite for efficient virtualisation
  - requires VM ISA (almost) same as processor ISA

Guest

```
ld    r0, curr_thrd
ld    r1, (r0,ASID)
mv    CPU_ASID, r1
ld    sp, (r1,kern_stk)
```

Exception

Hypervisor

```
lda   r1, vm_reg_ctxt
ld    r2, (r1,ofs_r0)
sto   r2, (r1,ofs_ASID)
```

UNSW
THE UNIVERSITY OF NEW SOUTH WALES

# Trap-and-Emulate Requirements
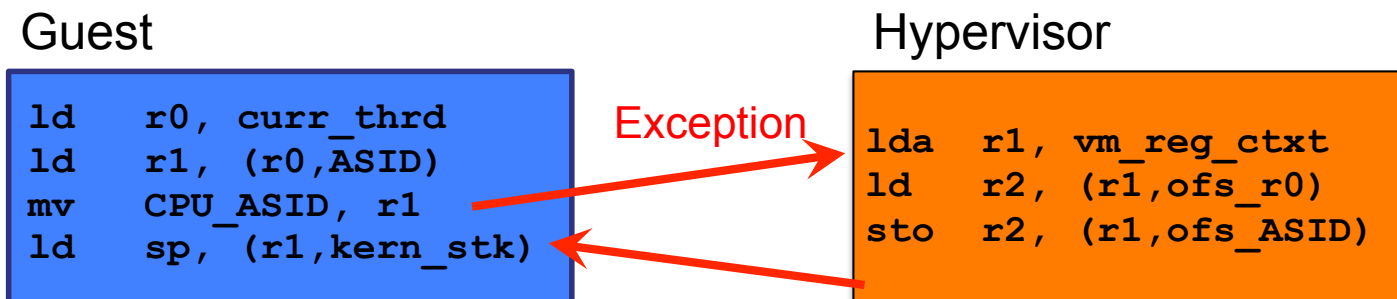
**Definitions:**

- **Privileged instruction:** traps when executed in user mode
  - Note: NO-OP is insufficient!
- **Privileged state:** determines resource allocation
  - Includes privilege mode, addressing context, exception vectors…
- **Sensitive instruction:** control- or behaviour-sensitive
  - **control sensitive:** changes privileged state
  - **behaviour sensitive:** exposes privileged state
    - incl instructions which are NO-OPs in user but not privileged state
- **Innocuous instruction:** not sensitive

- Some instructions are inherently sensitive
  - eg TLB load
- Others are context-dependent
  - eg store to page table
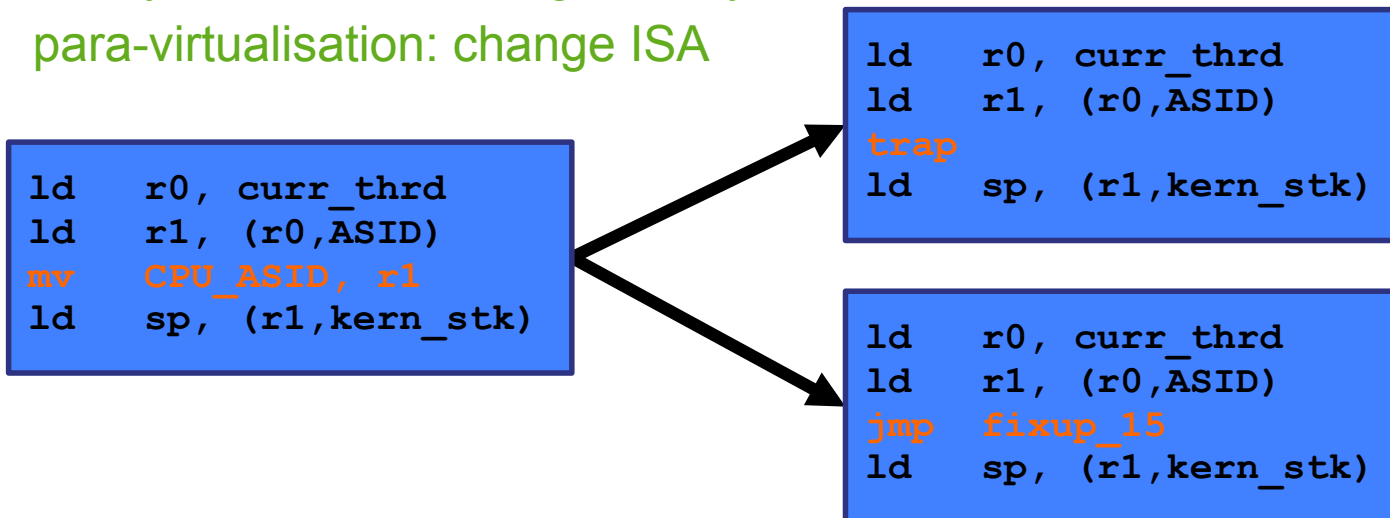
# Trap-and-Emulate Architectural Requirements

- ***T&E virtualisable*:** all sensitive instructions are privileged
  - Can achieve accurate, efficient guest execution
    - … by simply running guest binary on hypervisor
  - VMM controls resources
  - Virtualized execution indistinguishable from native, except:
    - resources more limited (smaller machine)
    - timing differences (if there is access to real time clock)

- **Recursively virtualisable**:
  - run hypervsior in VM
  - possible if hypervsior not timing dependent

Guest

```
ld    r0, curr_thrd
ld    r1, (r0,ASID)
mv    CPU_ASID, r1
ld    sp, (r1,kern_stk)
```

Exception

Hypervisor

```
lda   r1, vm_reg_ctxt
ld    r2, (r1,ofs_r0)
sto   r2, (r1,ofs_ASID)
```

UNSW
THE UNIVERSITY OF NEW SOUTH WALES

# Impure Virtualization

- Virtualise other than by T&E of unmodified binary
- Two reasons:
  - Architecture not T&E virtualisable
  - Reduce virtualisation overheads
- Change guest OS, replacing sensitive instructions
  - by trapping code (hypercalls)
  - by in-line emulation code
- Two approaches
  - binary translation: change binary
  - para-virtualisation: change ISA

```
ld    r0, curr_thrd
ld    r1, (r0,ASID)
mv    CPU_ASID, r1
ld    sp, (r1,kern_stk)
```

```
ld    r0, curr_thrd
ld    r1, (r0,ASID)
trap
ld    sp, (r1,kern_stk)
```

```
ld    r0, curr_thrd
ld    r1, (r0,ASID)
jmp   fixup_15
ld    sp, (r1,kern_stk)
```
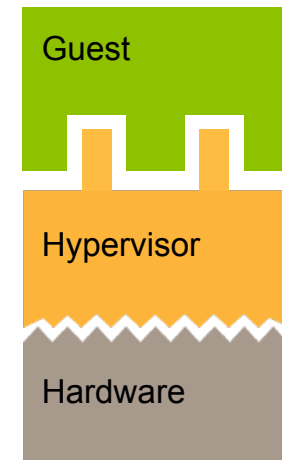
# Binary Translation

NICTA

- Locate sensitive instructions in guest binary,
  replace on-the-fly by emulation or trap/hypercall
  - pioneered by VMware
  - detect/replace combination of sensitive instruction for performance
  - modifies binary at load time, no source access required
- Looks like pure virtualisation!
- Very tricky to get right (especially on x86!)
  - Assumptions needed about sane guest behaviour
  - "Heroic effort" [Orran Krieger, then IBM, later VMware] 😊

UNSW
THE UNIVERSITY OF NEW SOUTH WALES

# Para-Virtualization

- New(ish) name, old technique
    - coined by Denali [Whitaker '02], popularised by Xen [Barham '03]
    - Mach Unix server [Golub '90], L4Linux [Härtig '97], Disco [Bugnion '97]
- Idea: manually port guest OS to modified (more high-level) ISA
    - Augmented by explicit hypervisor calls (hypercalls)
        - higher-level ISA to reduce number of traps
        - remove unvirtualisable instructions
        - remove "messy" ISA features which complicate
    - Generally outperforms pure virtualisation, binary re-writing
- Drawbacks:
    - Significant engineering effort
    - Needs to be repeated for each guest-ISA-hypervisor combination
    - Para-virtualised guests must be kept in sync with native evolution
    - Requires source

Guest

Hypervisor

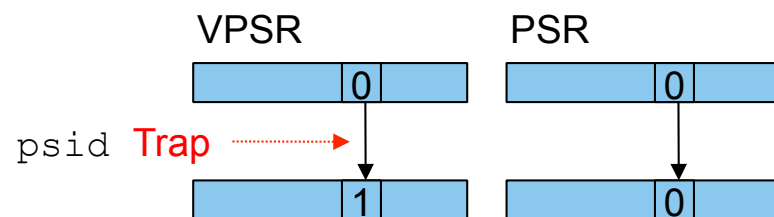Hardware

# Virtualization Overheads

NICTA

- VMM must maintain virtualised privileged machine state
  - processor status
  - addressing context
  - device state
- VMM needs to emulate privileged instructions
  - translate between virtual and real privileged state
  - eg guest ↔ real page tables
- Virtualisation traps are expensive
  - >1000 cycles on some Intel processors!
- Some OS operations involve frequent traps
  - STI/CLI for mutual exclusion
  - frequent page table updates during fork()
  - MIPS KSEG addresses used for physical addressing in kernel

UNSW
THE UNIVERSITY OF NEW SOUTH WALES

# Virtualization Techniques

- Impure virtualisation methods enable new optimisations
  - due to ability to control the ISA

- Example: Maintain some virtual machine state inside the VM
  - eg interrupt-enable bit (in virtual PSR)
  - requires changing guest's idea of where this bit lives
  - hypervisor knows about VM-local virtual state
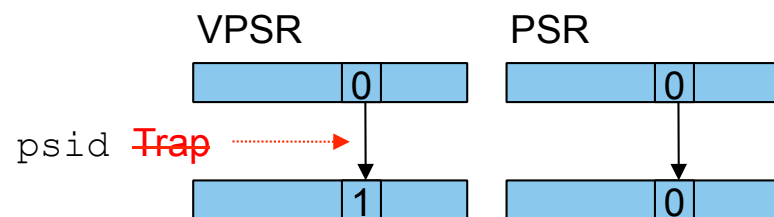    - eg queue vitual interrupt until guest enables in virtual PSR

VPSR                    PSR

| | 0 | |            | | 0 | |

psid Trap ---->

| | 1 | |            | | 0 | |

```
mov  r1,#VPSR
ldr  r0,[r1]
orr  r0,r0,#VPSR_ID
sto  r0,[r1]
```

# Virtualization Techniques

- Example: Lazy update of virtual machine state
  - virtual state is kept inside hypervisor
  - shadowed by copy inside VM
  - allow temporary inconsistency between primary and shadow
  - synchronise on next forced hypervsior invocation
    - actual trap
    - explicity hypercall when physical state must be updated
  - Example: guest enables FPU, handled lazily by hypervisor:
    - guest sets virtual FPU-enable bit
    - hypervisor synchronises on virtual kernel exit
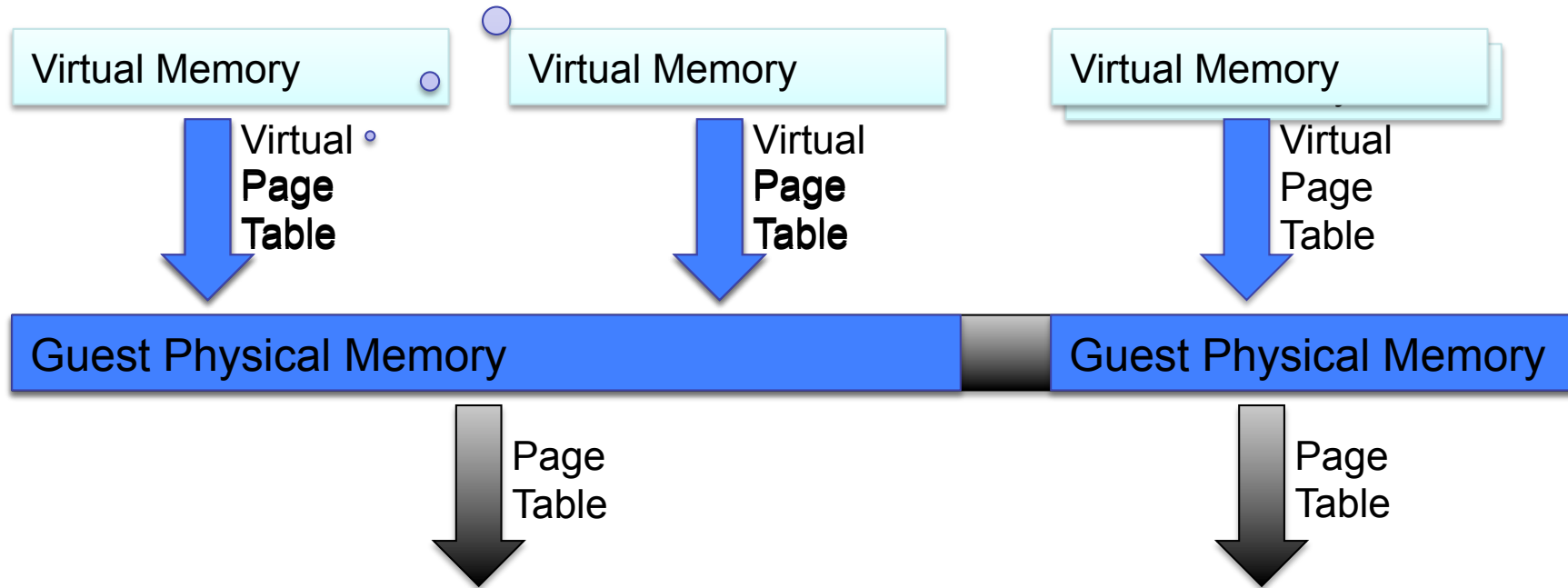- More examples later

VPSR            PSR

| 0 |          | 0 |

psid  ~~Trap~~ ┈┈►

| 1 |          | 0 |

```
mov  r1,#VPSR
ldr  r0,[r1]
orr  r0,r0,#VPSR_ID
sto  r0,[r1]
```

# Virtualization and Address Translation

Two levels of address translation!

| Virtual Memory | Virtual Memory | Virtual Memory |
|---|---|---|
| Virtual **Page Table** | Virtual **Page Table** | Virtual Page Table |

Guest Physical Memory | Guest Physical Memory

Page Table

Page Table

**Must implement with single MMU translation!**

# Virtualization Mechanics: Shadow Page Table

User
`ld r0, adr`

Guest virtual address

(Virtual) guest page table

Shadow (real) guest page table, translations cached in TLB

Hypervisor's guest memory map

Virt PT ptr (Software)

Guest OS

Guest physical address

PT ptr (Hardware)

Hypervisor

Physical address

Memory

data

NICTA

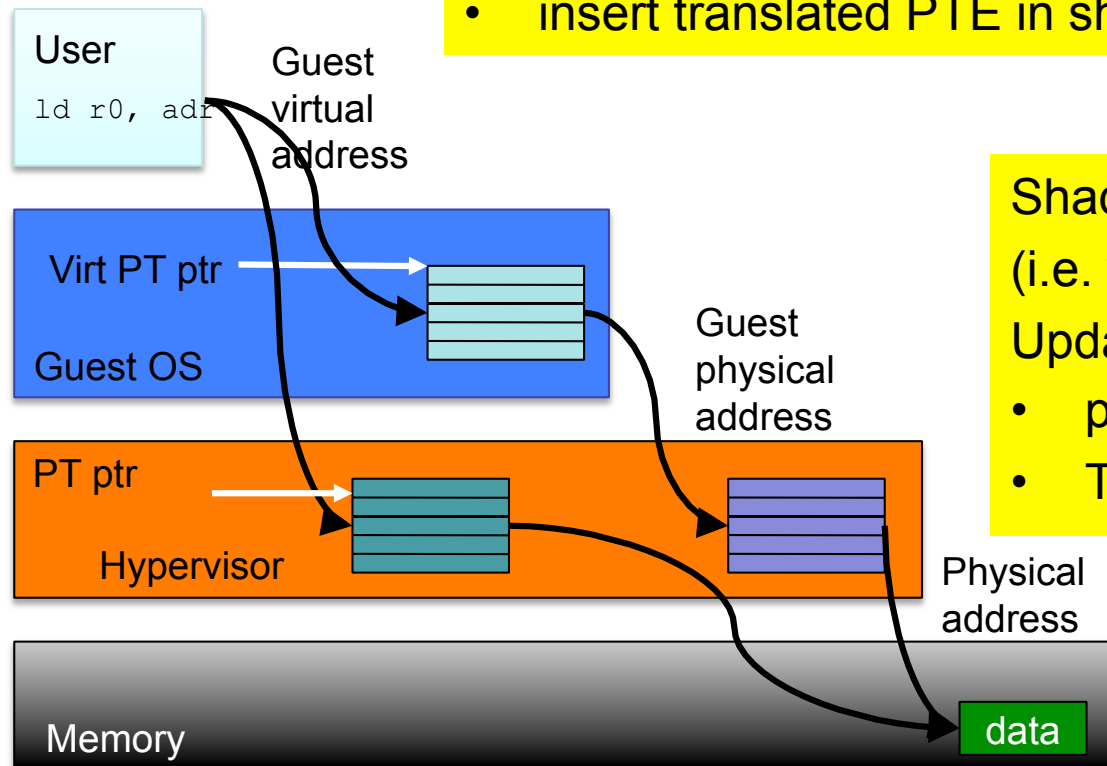UNSW
THE UNIVERSITY OF NEW SOUTH WALES

# Virtualization Mechanics: Shadow Page Table

**NICTA**

Hypervisor must shadow (virtualize)
all PT updates by guest:

- trap guest writes to guest PT
- translate guest PA in guest (virtual) PTE using guest memory map
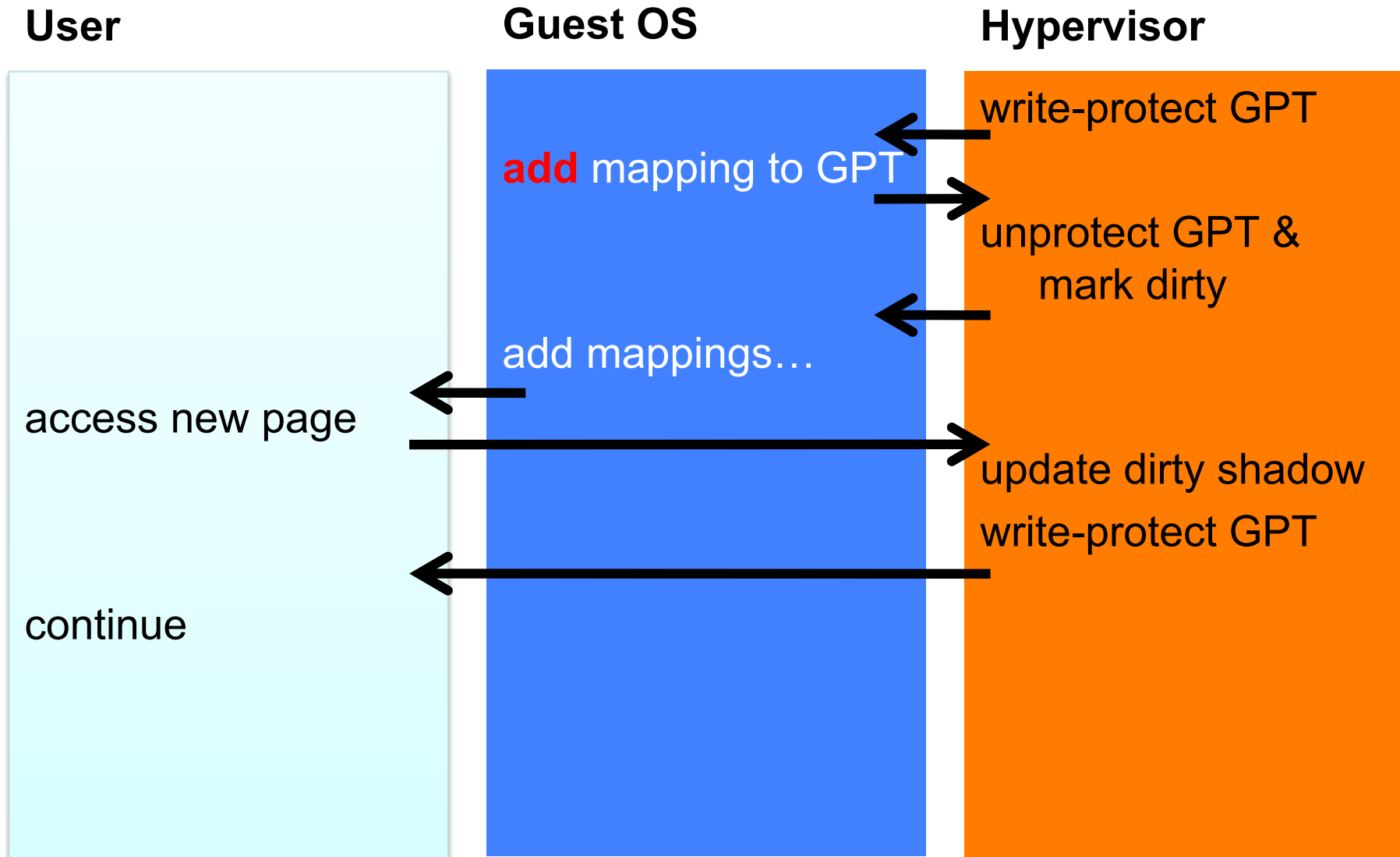- insert translated PTE in shadow PT

Used by VMware
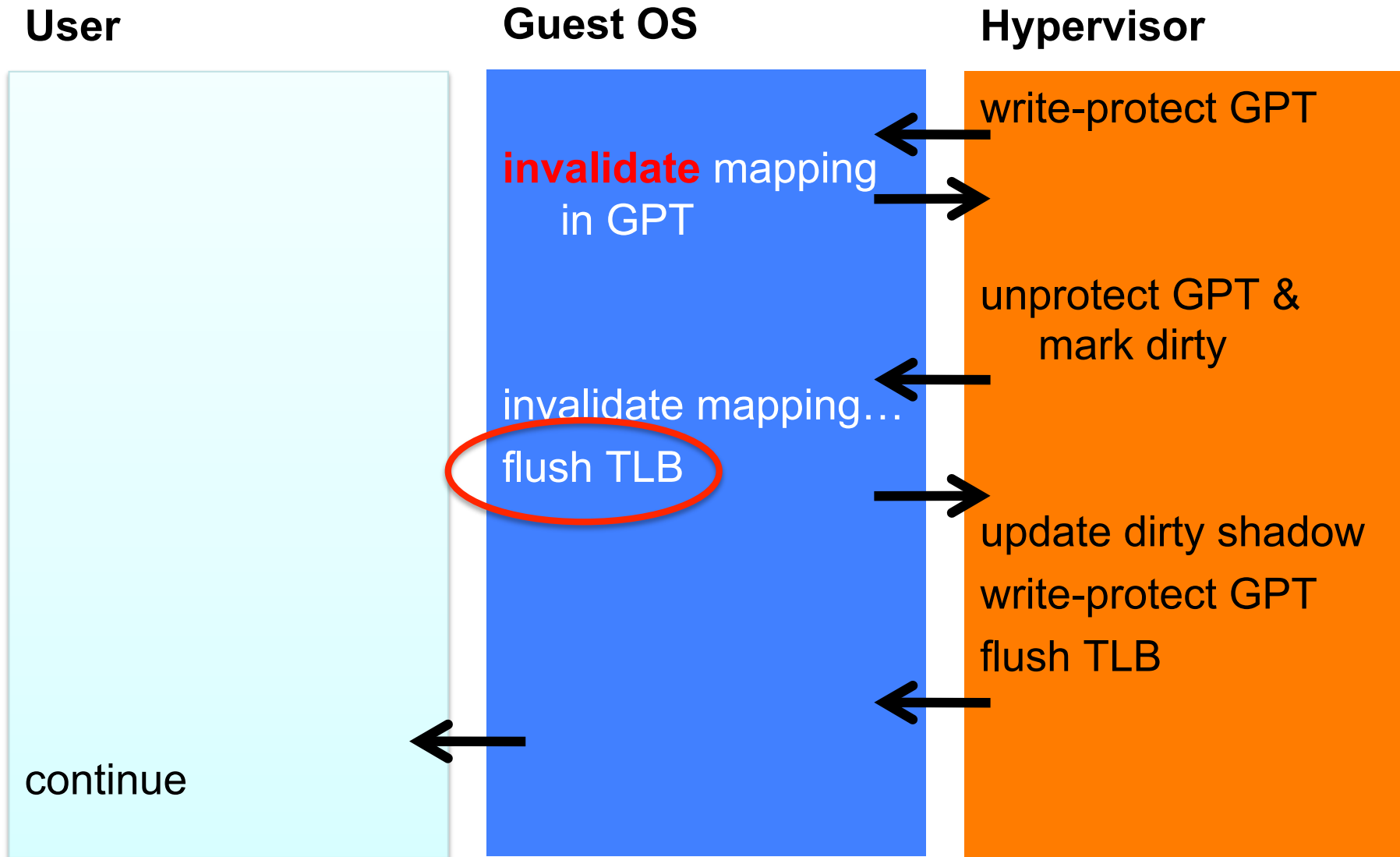
Shadow PT has TLB semantics
(i.e. weak consistency) ⇒
Update at synchronisation points:

- page faults
- TLB flushes

User
ld r0, adr

Guest virtual address

Virt PT ptr

Guest OS

Guest physical address

PT ptr

Hypervisor

Physical address

Memory

data

**UNSW**
THE UNIVERSITY OF NEW SOUTH WALES

# Virtualisation Semantics: Lazy Shadow Update

**User**        **Guest OS**        **Hypervisor**

write-protect GPT

**add** mapping to GPT

unprotect GPT &
mark dirty

add mappings…

access new page

update dirty shadow
write-protect GPT

continue

   

# Virtualisation Semantics: Lazy Shadow Update

**User**

**Guest OS**

**Hypervisor**

write-protect GPT

**invalidate** mapping
in GPT

unprotect GPT &
mark dirty

invalidate mapping…

flush TLB

update dirty shadow

write-protect GPT

flush TLB

continue

UNSW
THE UNIVERSITY OF NEW SOUTH WALES

# Virtualization Mechanics: Real Guest PT

**Hypervisor maintains guest PT**

User
`ld r0, adr`

Guest virtual address

Guest OS

Hypervisor

Guest PT

HV PT

Physical address

Memory

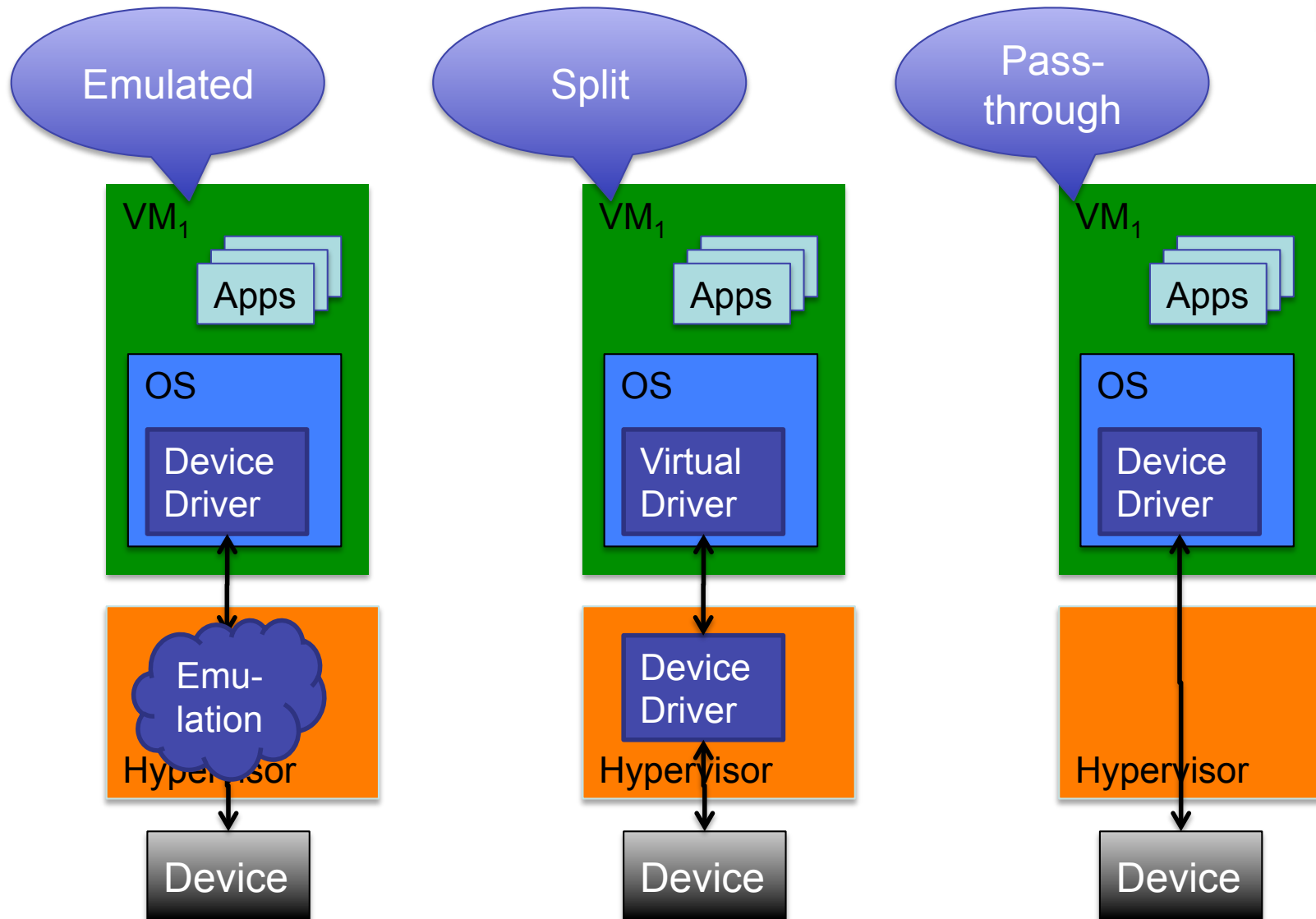data

- On guest PT access must translate (virtualize) PTEs
  - store: translate guest "PTE" to real PTE
  - load: translate real PTE to guest "PTE"
- Each guest PT access traps!
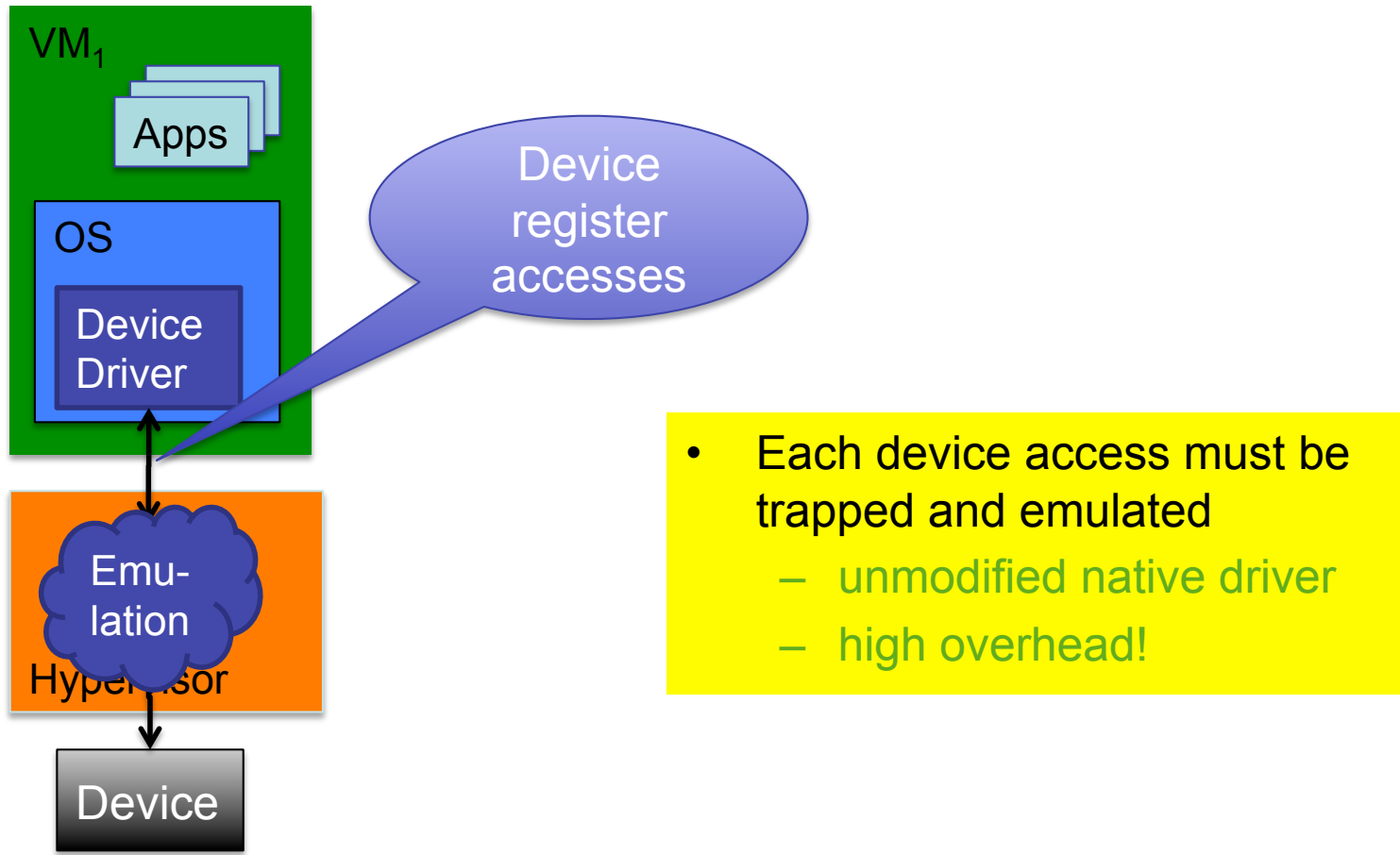  - including reads
  - high overhead

NICTA

UNSW
THE UNIVERSITY OF NEW SOUTH WALES

# Virtualization Mechanics: Optimised Guest PT

**NICTA**

Para-virtualized guest "knows" it is virtualized

User
`ld r0, adr`

Guest virtual address

Guest OS

Hypervisor

Guest PT        HV PT

Physical address

Memory     data

- Guest translates PTEs itself when reading from PT
  - supported by Linux PT-access wrappers
- Guest batches PT updates using hypercalls
  - reduced overhead

Used by original Xen

**UNSW**
THE UNIVERSITY OF NEW SOUTH WALES

# Virtualization Mechanics: 3 Device Models



© 2012 Gernot Heiser UNSW/NICTA. Distributed under Creative Commons Attribution License

# Virtualization Mechanics: Emulated Device

VM$_1$

Apps

OS

Device Driver

Device register accesses

Emu-lation

Hypervisor

Device

- Each device access must be trapped and emulated
  - unmodified native driver
  - high overhead!

UNSW
THE UNIVERSITY OF NEW SOUTH WALES

# Virtualization Mechanics: Split Driver (Xen speak)



- Simplified, high-level device interface
  - small number of hypercalls
  - new (but very simple) driver
  - low overhead
  - must port drivers to hypervisor

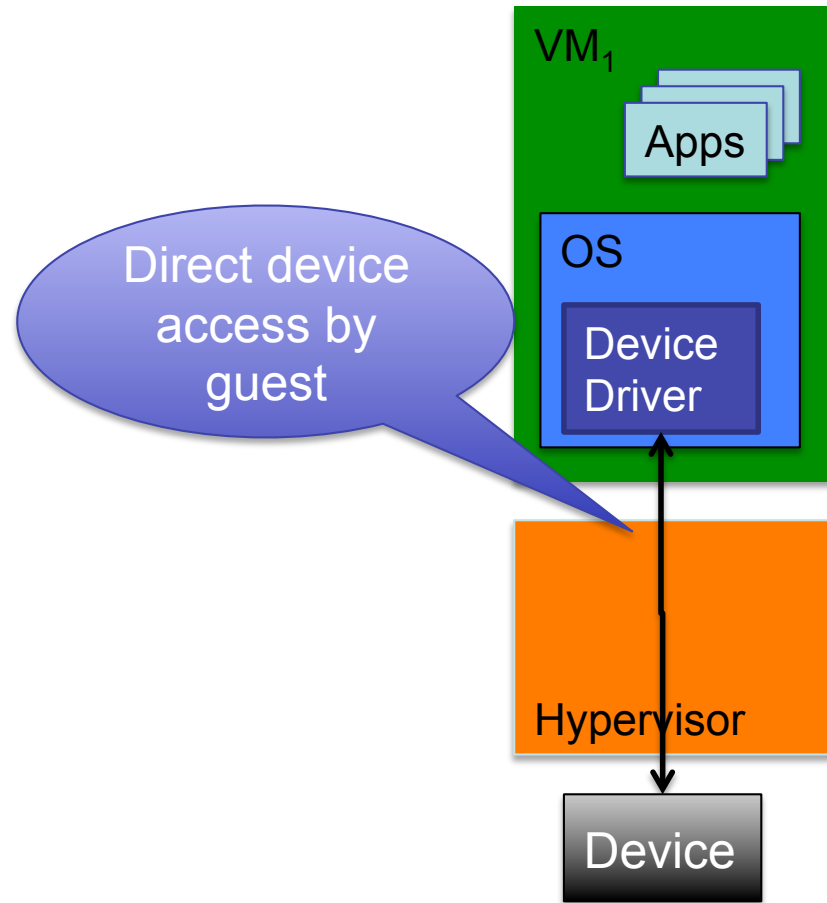# Virtualization Mechanics: Driver OS (Xen Dom0)



- **Leverage Driver-OS native drivers**
  - no driver porting
  - must trust complete Driver OS guest!

# Virtualization Mechanics: Pass-Through Driver

- **Unmodified native driver**

- **Must trust driver (and guest)**
  - unless have hardware support (I/O MMU)

VM$_1$

Apps

OS

Device Driver

Direct device access by guest

Hypervisor

Device

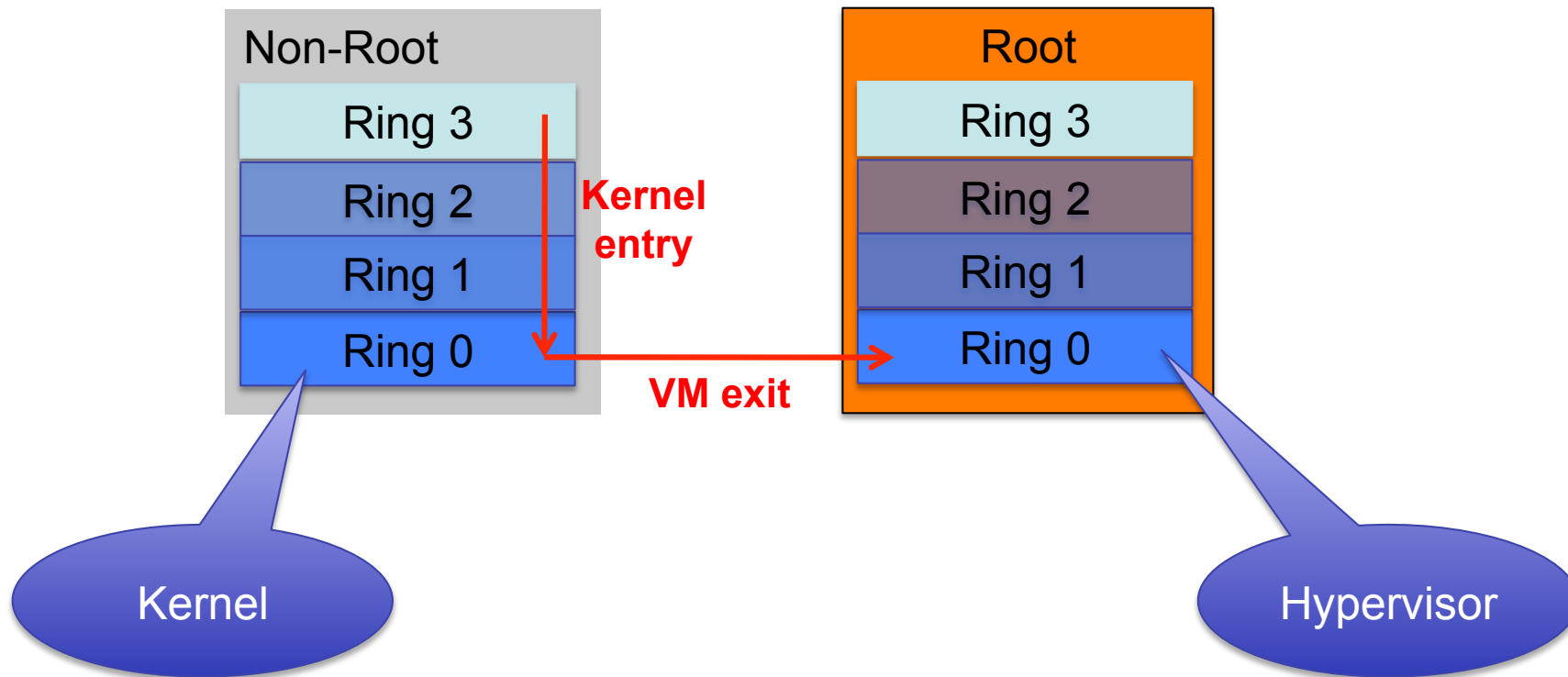# Modern Architectures Not T&E Virtualisable

- Examples:
  - x86: many non-virtualizable features
    - e.g. sensitive PUSH of PSW is not privileged
    - segment and interrupt descriptor tables in virtual memory
    - segment description expose privileged level
  - MIPS: mostly ok, but
    - kernel registers k0, k1 (for save/restore state) user-accessible
    - performance issue with virtualising KSEG addresses
  - ARM: mostly ok, but
    - some instructions undefined in user mode (banked registers, CPSR)
    - PC is a GPR, exception return in MOVS to PC, doesn't trap
- Addressed by virtualization extensions to ISA
  - x86, Itanium since ~2006 (VT-x, VT-i), ARM since '12
  - additional processor modes and other features
  - all sensitive ops trap into hypervisor or made innocuous (shadow state)
    - eg guest copy of PSW

# x86 Virtualization Extensions (VT-x)

- New processor mode: *VT-x root mode*
  - orthogonal to protection rings
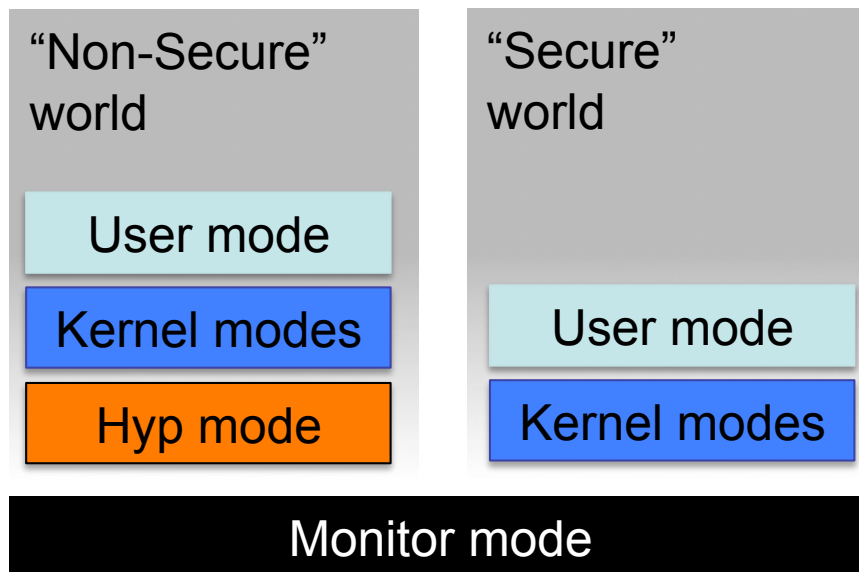  - entered on virtualisation trap
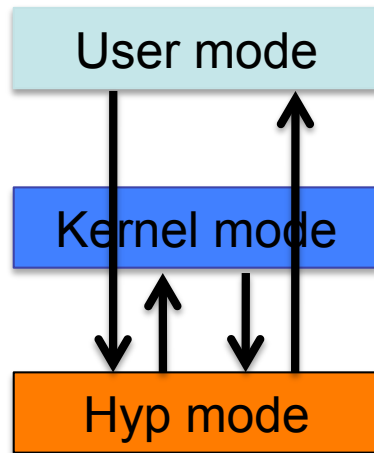
# ARM Virtualization Extensions (1)

## Hyp mode

- New privilege level
  - Strictly higher than kernel
  - Virtualizes or traps *all* sensitive instructions
  - Only available in ARM TrustZone "non-secure" mode
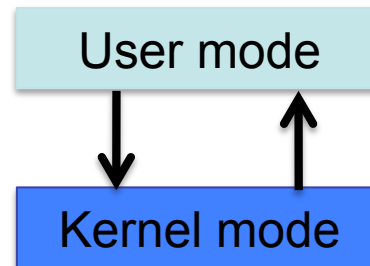
| "Non-Secure" world | "Secure" world |
|---|---|
| User mode | |
| Kernel modes | User mode |
| Hyp mode | Kernel modes |

Monitor mode

© 2012 Gernot Heiser UNSW/NICTA. Distributed under Creative Commons Attribution License

# ARM Virtualization Extensions (2)

## Configurable Traps

x86 similar

User mode

Kernel mode

*Native syscall*

Can configure traps to go directly to guest OS

User mode

Kernel mode

Hyp mode

*Virtual syscall*

User mode

Kernel mode

Hyp mode

*Virtual syscall*
*Trap to guest*

NICTA

UNSW
THE UNIVERSITY OF NEW SOUTH WALES

# ARM Virtualization Extensions (3)

## Emulation

1) Load faulting instruction
   - Compulsory L1-D miss!
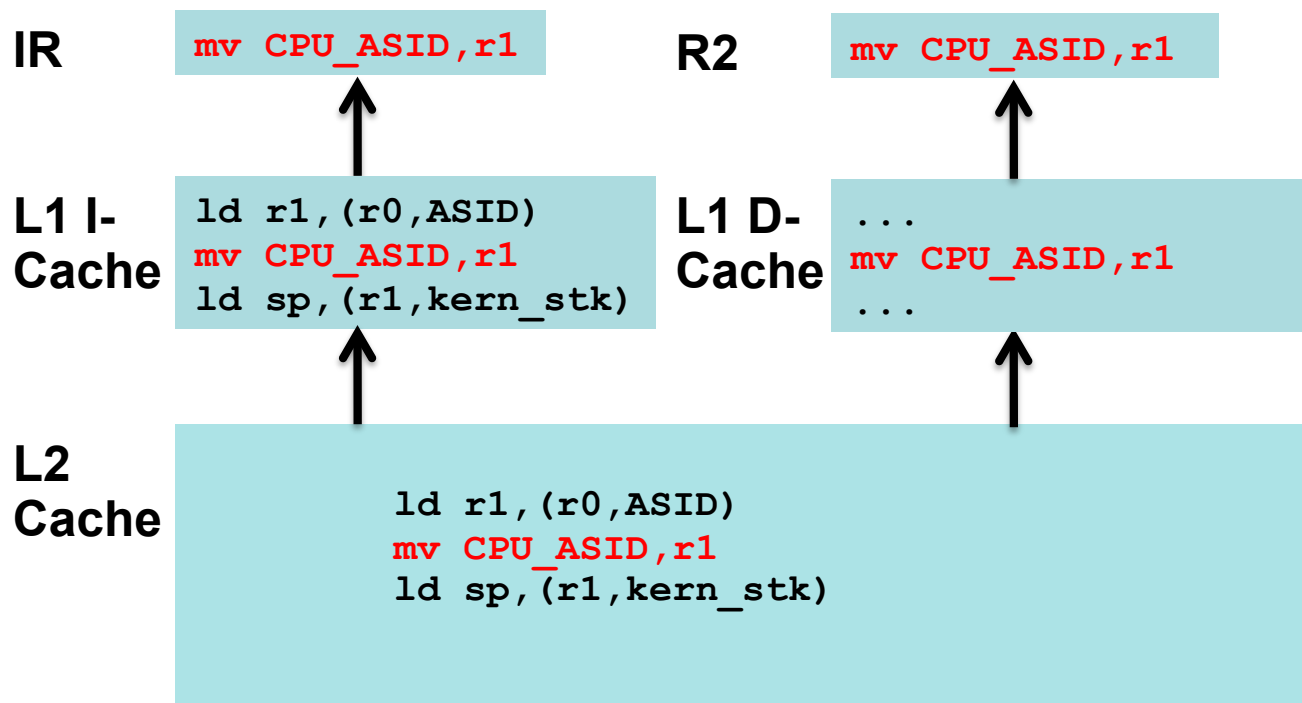2) Decode instruction
   - Complex logic
3) Emulate instruction
   - Usually straightforward

**IR**
```
mv CPU_ASID,r1
```

**R2**
```
mv CPU_ASID,r1
```

**L1 I-Cache**
```
ld r1,(r0,ASID)
mv CPU_ASID,r1
ld sp,(r1,kern_stk)
```

**L1 D-Cache**
```
...
mv CPU_ASID,r1
...
```

**L2 Cache**
```
ld r1,(r0,ASID)
mv CPU_ASID,r1
ld sp,(r1,kern_stk)
```

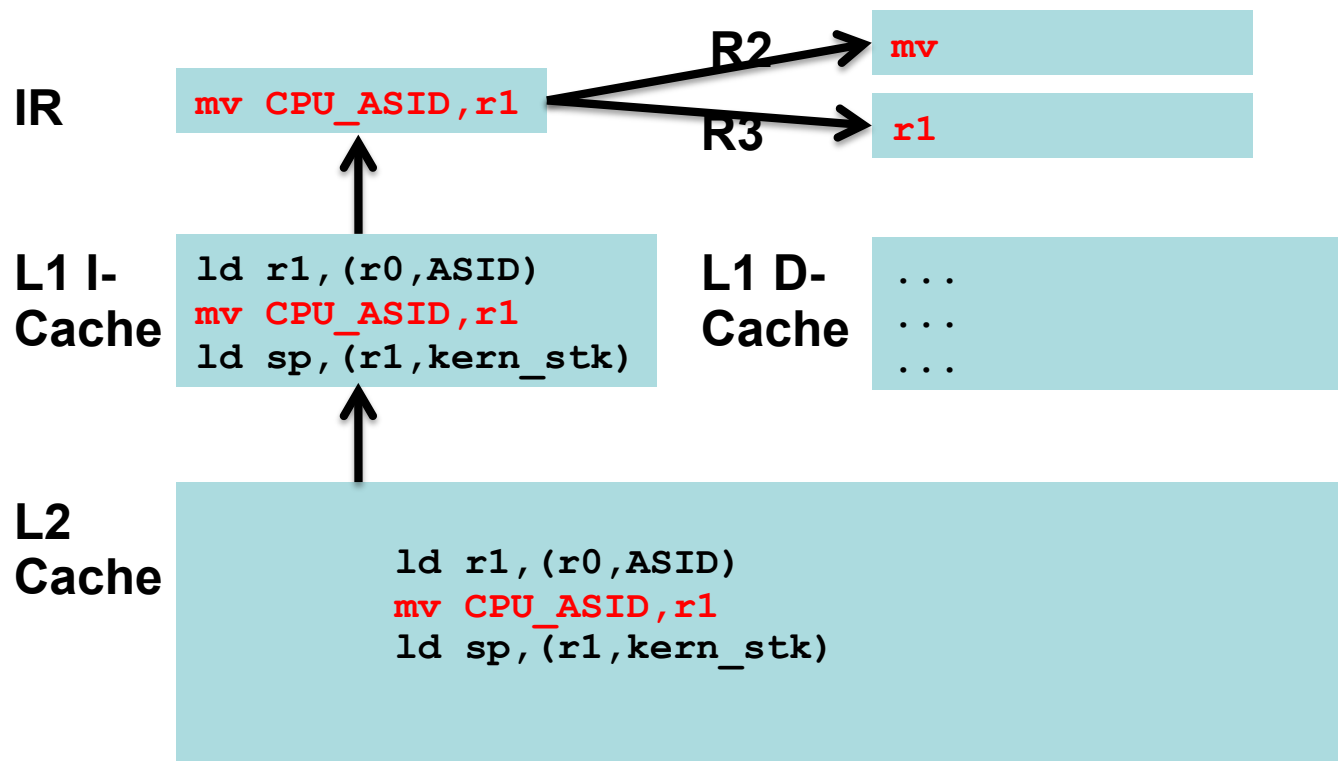UNSW
THE UNIVERSITY OF NEW SOUTH WALES

# ARM Virtualization Extensions (3)

**Emulation Support**

No x86 equivalent

– HW decodes instruction
  - No L1 miss
  - No software decode
– SW emulates instruction
  - Usually straightforward

**IR**

```
mv CPU_ASID,r1
```
R2 → `mv`
R3 → `r1`

**L1 I-Cache**

```
ld r1,(r0,ASID)
mv CPU_ASID,r1
ld sp,(r1,kern_stk)
```

**L1 D-Cache**

```
...
...
...
```

**L2 Cache**

```
ld r1,(r0,ASID)
mv CPU_ASID,r1
ld sp,(r1,kern_stk)
```

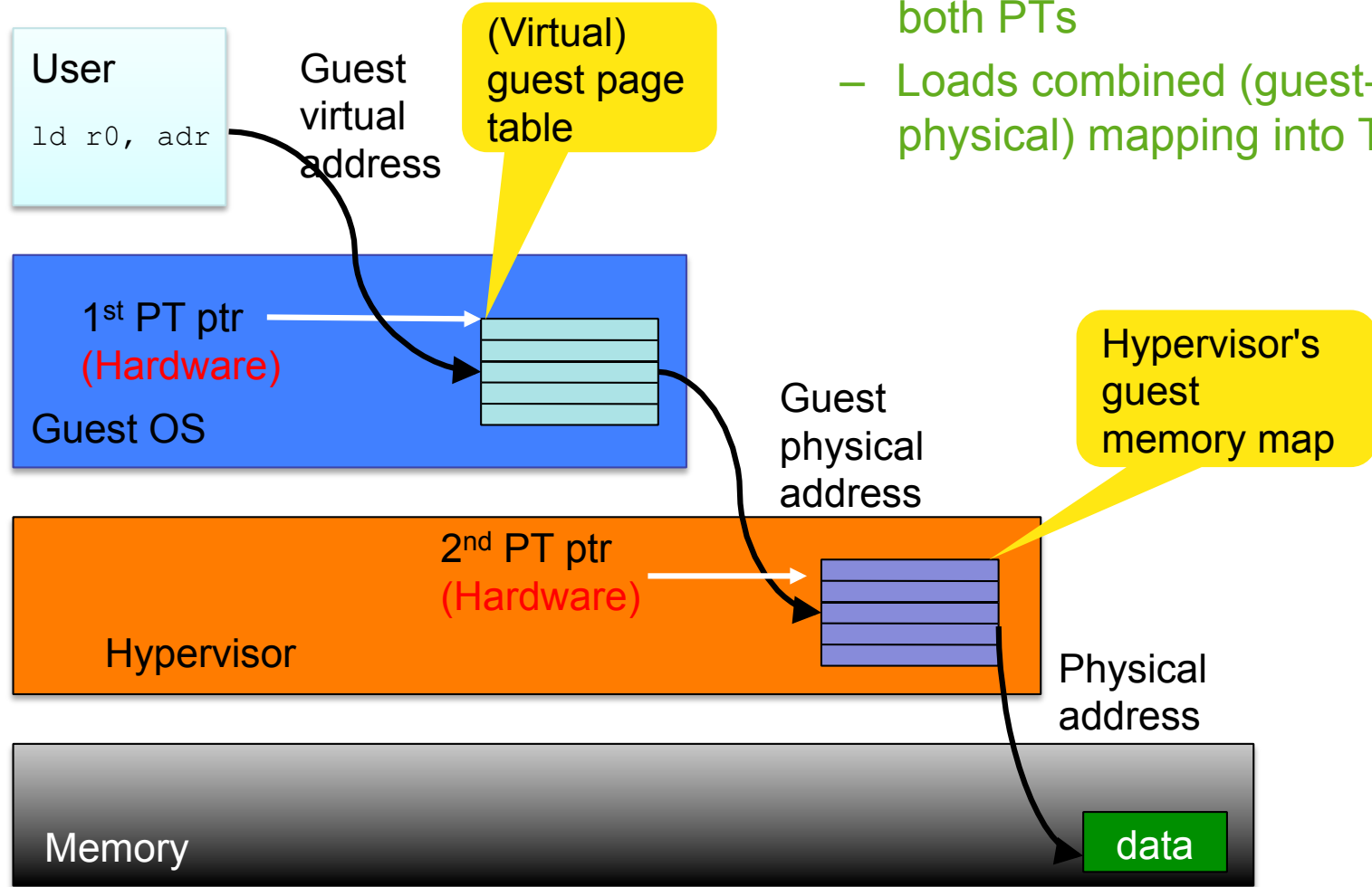UNSW
THE UNIVERSITY OF NEW SOUTH WALES

# ARM Virtualization Extensions (4)

## 2-stage translation

x86 similar

– Hardware PT walker traverses both PTs
– Loads combined (guest-virtual to physical) mapping into TLB

User
`ld r0, adr`

Guest virtual address

(Virtual) guest page table

1st PT ptr (Hardware)

Guest OS

Guest physical address

Hypervisor's guest memory map

2nd PT ptr (Hardware)

Hypervisor

Physical address

Memory

data

NICTA

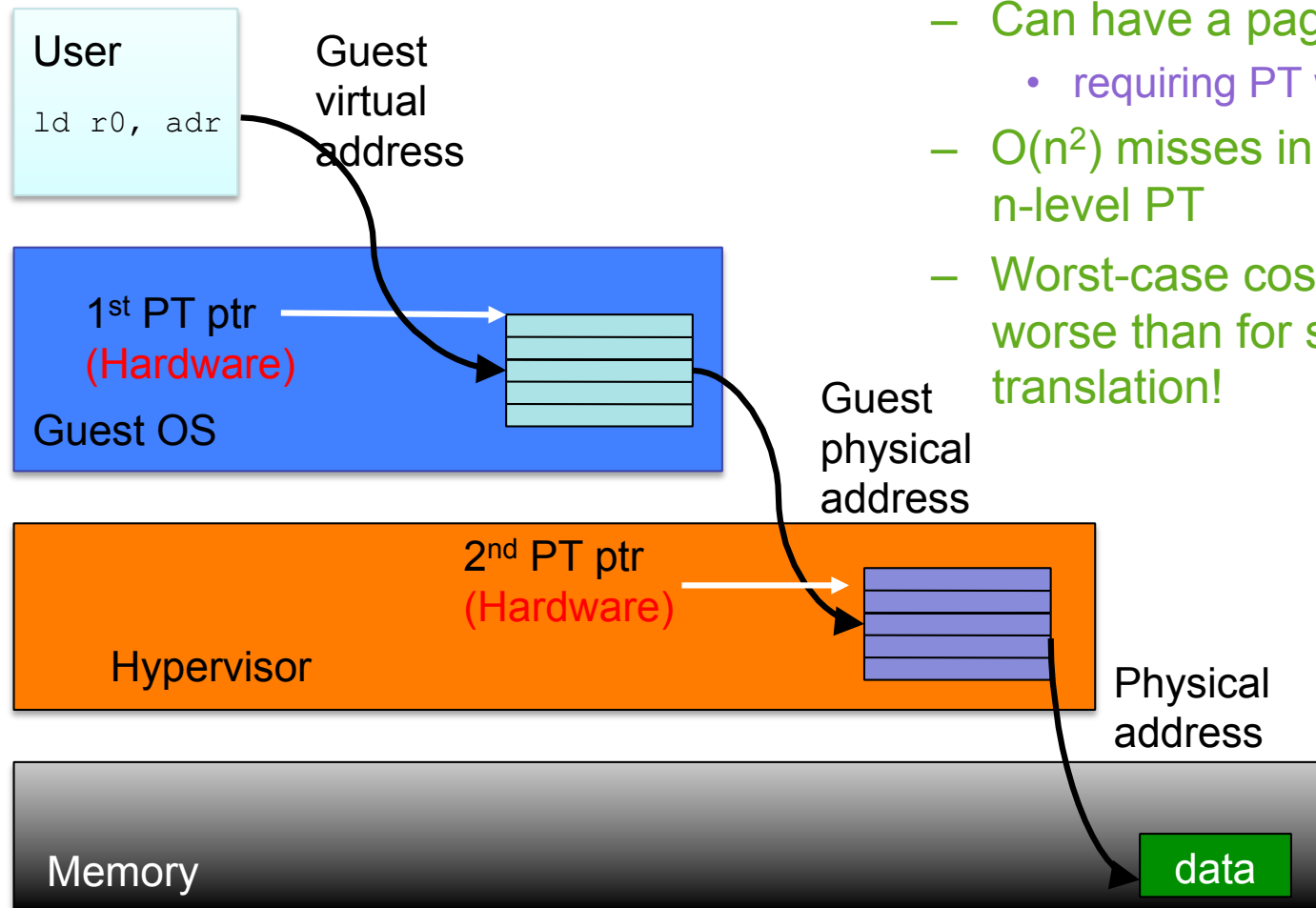UNSW
THE UNIVERSITY OF NEW SOUTH WALES

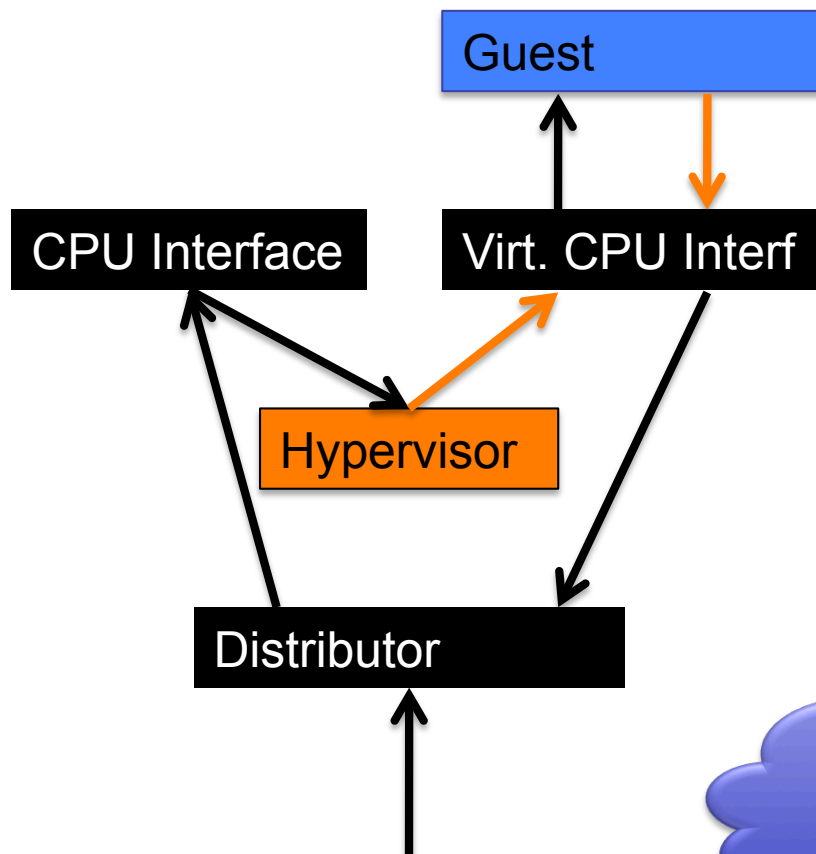# ARM Virtualization Extensions (4)

## 2-stage translation cost

– On page fault walk twice number of page tables!

– Can have a page miss on each
  • requiring PT walk

– $O(n^2)$ misses in worst case for n-level PT

– Worst-case cost is massively worse than for single-level translation!

User
`ld r0, adr`

Guest virtual address

1st PT ptr
(Hardware)

Guest OS

Guest physical address

2nd PT ptr
(Hardware)

Hypervisor

Physical address

Memory

data

UNSW
THE UNIVERSITY OF NEW SOUTH WALES

# ARM Virtualization Extensions (5)
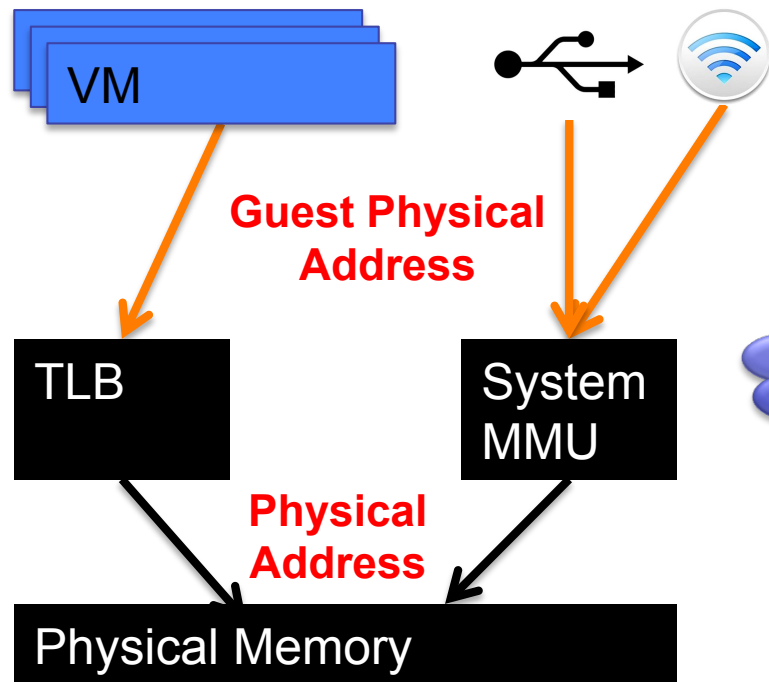
## Virtual Interrupts

- ARM has 2-part IRQ controller
  - Global "distributor"
  - Per-CPU "interface"
- New H/W "virt. CPU interface"
  - Mapped to guest
  - Used by HV to forward IRQ
  - Used by guest to acknowledge
- Halves hypervisor invocations for interrupt virtualization

```
                    Guest

CPU Interface      Virt. CPU Interf

        Hypervisor

            Distributor
```

x86: issue only for legacy level-triggered IRQs

© 2012 Gernot Heiser UNSW/NICTA. Distributed under Creative Commons Attribution License

# ARM Virtualization Extensions (6)

## System MMU (I/O MMU)



- Devices use virtual addresses
- Translated by system MMU
  - elsewhere called I/O MMU
  - translation cache, like TLB
  - reloaded from same page table

x86 different (VT-d)

- Can do pass-through I/O safely
  - guest accesses device registers
  - no hypervisor invocation

# Hypervisor Size

| Hypervisor | ISA | Type | Kernel | User |
|------------|-----|------|-------:|-----:|
| OKL4 | ARMv7 | para-virtualization | 9.8 kLOC | 0 |
| *Prototype* | ARMv7 | pure virtualization | 6 kLOC | 0 |
| Nova | x86 | pure virtualization | 9 kLOC | 27 kLOC |

- Size (& complexity) reduced about 40% wrt to para-virtualization
- Much smaller than x86 pure-virtualization hypervisor
  - Mostly due to greatly reduced need for instruction emulation

# Overheads (Estimated)

| Operation | Pure virtualization | | Para-virtualiz. |
| --- | --- | --- | --- |
| | Instruct | Cycles (est) | Cycles (approx) |
| Guest system call | 0 | 0 | 300 |
| Hypervisor entry + exit | 120 | 650 | 150 |
| IRQ entry + exit | 270 | 900 | 300–400? |
| Page fault | 356 | 1500 | 700 |
| Device emul. | 249 | 1040 | N/A |
| Device emul. (accel.) | 176 | 740 | N/A |
| World switch | 2824 | 7555 | 200 |

- No overhead on regular (virtual) syscall – unlike para-virtualization
- Invoking hypervisor 500–1200 cycles (0.6–1.5 µs) more than para
- World switch in ~10 µs compared to 0.25 µs for para
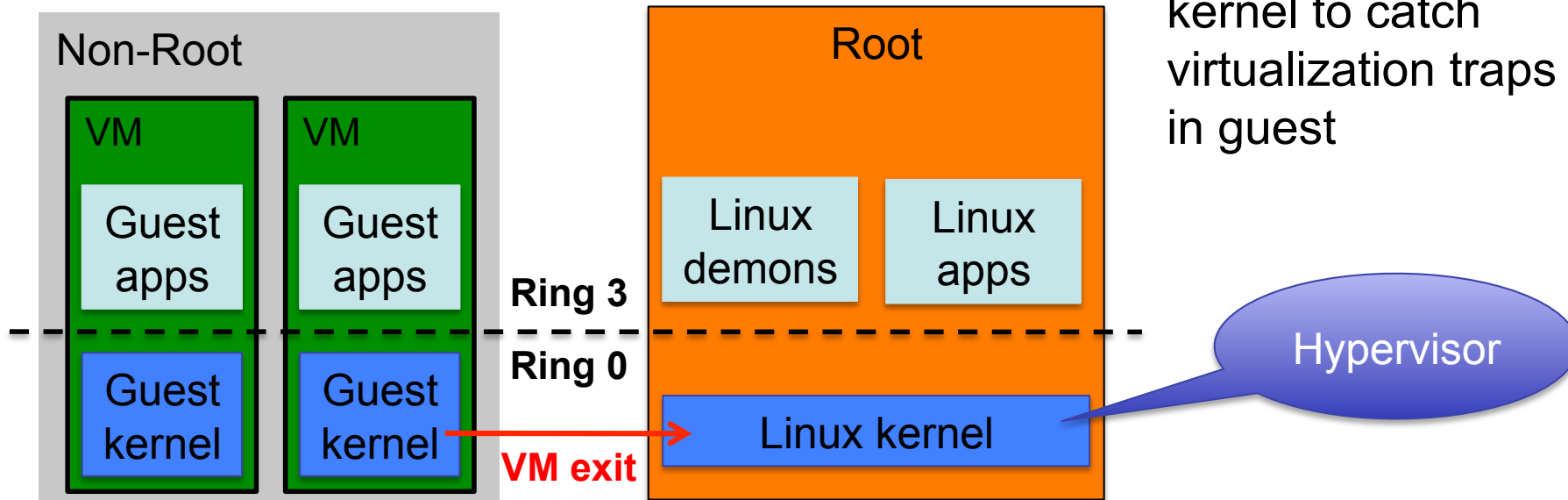- ⇒ Trade-offs differ

UNSW
THE UNIVERSITY OF NEW SOUTH WALES

# Hybrid Hypervisor OSes

- Idea: turn standard OS into hypervisor
  - … by running in VT-x root mode
  - eg: KVM ("kernel-based virtual machine")
- Can re-use Linux drivers etc
- Huge trusted computing base
- Often falsely called a Type-2 hypervisor

Variant: *VMware MVP*

- ARM hypervisor
  - pre-HW support
- re-writes exception vectors in Android kernel to catch virtualization traps in guest

**Non-Root**

| VM | VM |
|---|---|
| Guest apps | Guest apps |
| Guest kernel | Guest kernel |

**Root**

| Linux demons | Linux apps |
|---|---|
| Linux kernel | |

Ring 3

Ring 0

**VM exit**

Hypervisor

NICTA

UNSW
THE UNIVERSITY OF NEW SOUTH WALES

# Fun and Games with Hypervisors

- Time-travelling virtual machines [King '05]
  - debug backwards by replay VM from checkpoint, log state changes
- SecVisor: kernel integrity by virtualisation [Seshadri '07]
  - controls modifications to kernel (guest) memory
- Overshadow: protect apps from OS [Chen '08]
  - make user memory opaque to OS by transparently encrypting
- Turtles: Recursive virtualisation [Ben-Yehuda '10]
  - virtualize VT-x to run hypervisor in VM
- CloudVisor: mini-hypervisor underneath Xen [Zhang '11]
  - isolates co-hosted VMs belonging to different users
  - leverages remote attestation (TPM) and Turtles ideas

… and many more!

NICTA

UNSW
THE UNIVERSITY OF NEW SOUTH WALES

# Hypervisors vs Microkernels

- Both contain all code executing at highest privilege level
  - Although hypervisor may contain user-mode code as well
    - privileged part usually called "hypervisor"
    - user-mode part often called "VMM"
- Both need to abstract hardware resources
  - Hypervisor: abstraction closely models hardware
  - Microkernel: abstraction designed to support wide range of systems
- What must be abstracted?
  - Memory
  - CPU
  - I/O
  - Communication

Difference to traditional terminology!

NICTA

UNSW
THE UNIVERSITY OF NEW SOUTH WALES

# What's the difference?

| Resource | Hypervisor | Microkernel |
|---|---|---|
| Memory | Virtual MMU (vMMU) | Address space |
| CPU | Virtual CPU (vCPU) | Thread or scheduler activation |
| I/O | • Simplified virtual device<br>• Driver in hypervisor<br>• Virtual IRQ (vIRQ) | • IPC interface to user-mode driver<br>• Interrupt IPC |
| Communication | Virtual NIC, with driver and network stack | High-performance message-passing IPC |

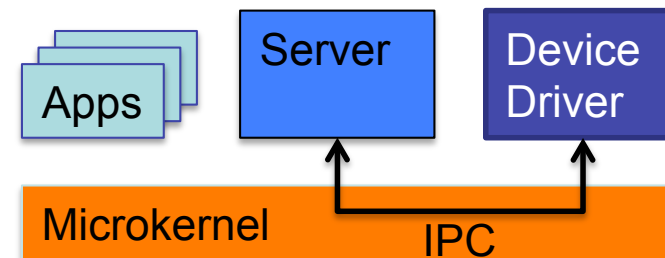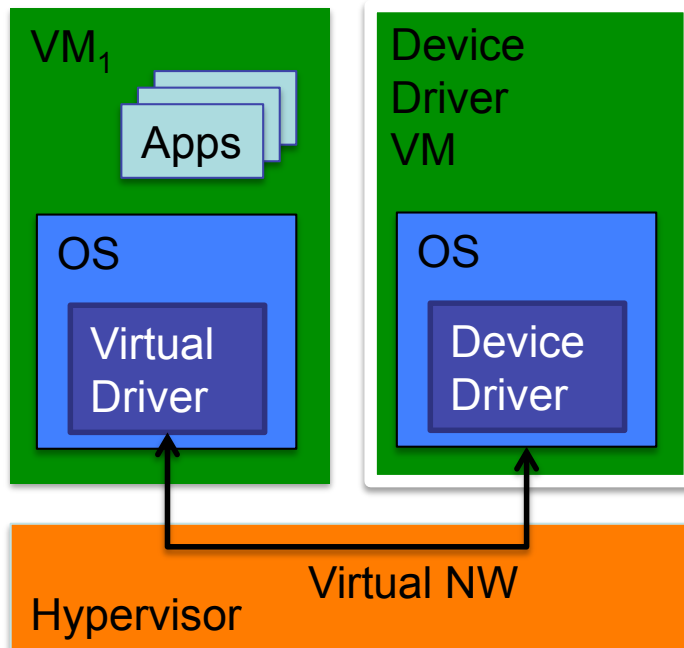Just page tables in disguise

Just kernel-scheduled activities

Real Difference?

- Similar abstractions
- Optimised for different use cases

Modelled on HW, Re-uses SW

Minimal overhead, Custom API

NICTA

UNSW
THE UNIVERSITY OF NEW SOUTH WALES

# Closer Look at I/O and Communication



- Communication is critical for I/O
  - Microkernel IPC is highly optimised
  - Hypervisor inter-VM communication is frequently a bottleneck

# Hypervisors vs Microkernels: Drawbacks

**NICTA**

## Hypervisors:

- Communication is Achilles heel
  - more important than expected
    - critical for I/O
  - plenty improvement attempts in Xen


- Most hypervisors have big TCBs
  - infeasible to achieve high assurance of security/safety
  - in contrast, microkernel implementations can be proved correct

## Microkernels:

- Not ideal for virtualization
  - API not very effective
    - L4 virtualization performance close to hypervisor
    - effort much higher
  - Virtualization needed for legacy

- L4 model uses kernel-scheduled threads for more than exploiting parallelism
  - Kernel imposes policy
  - Alternatives exist, eg. K42 uses scheduler activations

**UNSW**
THE UNIVERSITY OF NEW SOUTH WALES