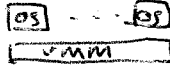


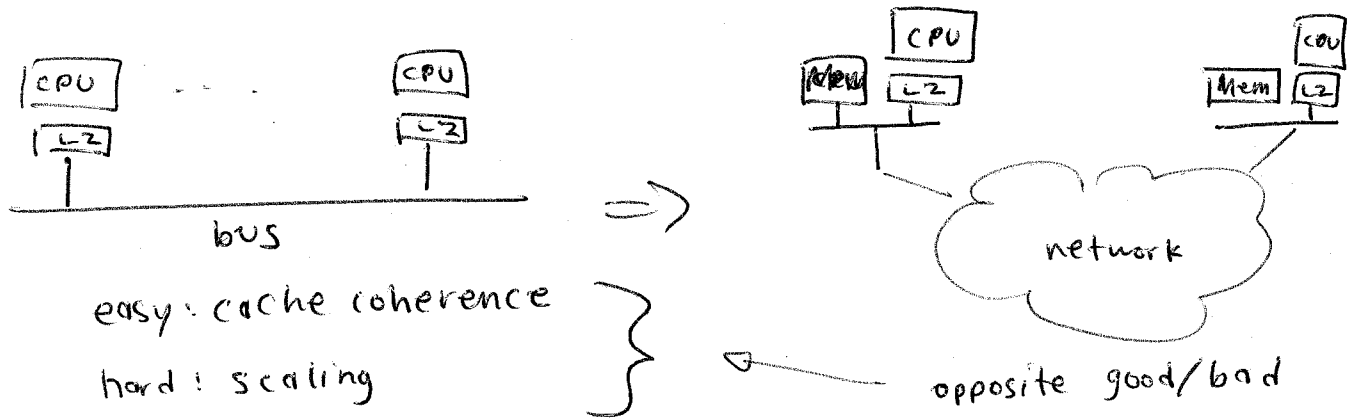
Disco

(Spr '05)

Overview

- Beautiful system, hacking, + paper!
- Illusions, Layering \Rightarrow  or deal w/ NUMA here
- Problems: Overhead, Sharing, Resource Mgmt

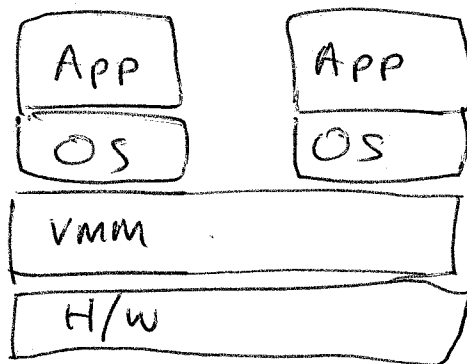
Background: ccNUMA



Question: How to build OS for ccNUMA?

\Rightarrow Hard: Data structures in OS, page placement, etc.
OS is big/ugly/hard to change

Solution: Virtual Machines



\Rightarrow Advantages?

- Hide tough issues ("parallel") from OS
- Can run different OS's concurrently
- Portability layer

Problems:

Overhead
Sharing
Resource Mgmt
(Breaking the Illusion)

Why Overheads?

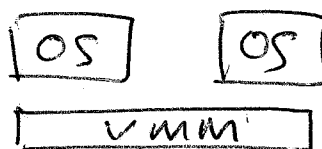
- Time

- Space



could be a simulator!

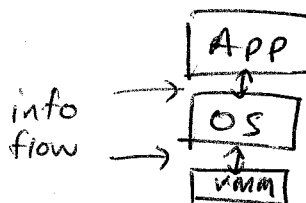
(but usually can just directly execute instrs.)



{ file cache
page tables
etc. }

Resource Mgmt Problems

- Information flow



e.g. OS in idle loop,
page on free list

Sharing problems

- Old VM/370 → no file sharing allowed
- Here : Distributed Systems Technology

Disco Virtualization

3

CPU: MIPS R10K

"direct execution"

set real registers to VCPU registers, jump to VPC

hard: privileged instructions (TLB, phys mem, I/O)

MIPS details:

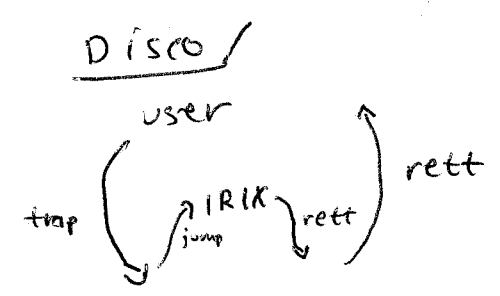
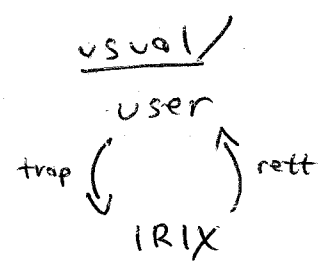
why?
MIPS wanted to support "modern" OS's
→ supervisor
kernel

usual/
App
(not used)
IRIX

Disco/
App
OS
Disco

Example: ~~3~~

System call:
trap/rett



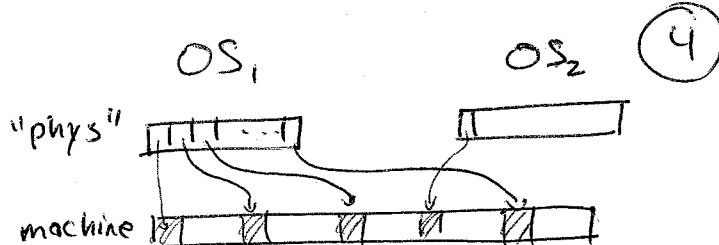
key: supervisor mode in MIPS

can access more mem than user mode,
but not privileged inst, phys mem

Result: new protection boundaries
in system

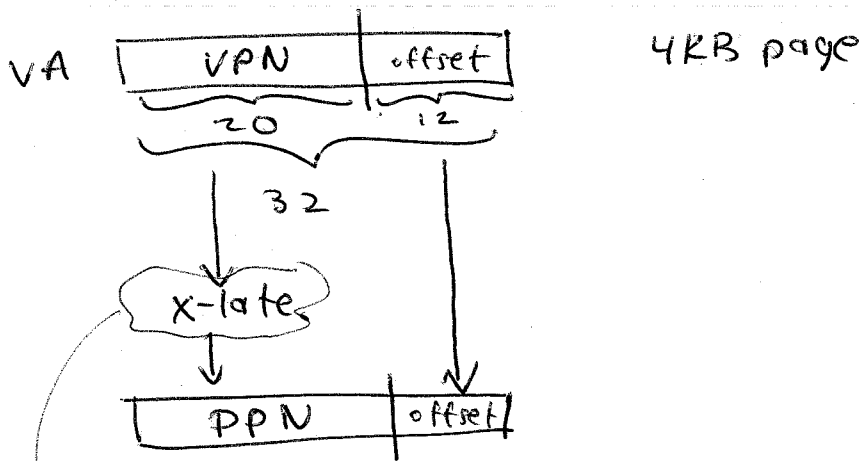
Memory

$VA \rightarrow PA \rightarrow MA$
usual OS w/ Disco



All done thru TLB ; e.g.:

\Rightarrow load $O(R_1), R_2$



Q: where are all x-lations held?
 A: page table (just a data structure in the OS)
 (VA \rightarrow PA)

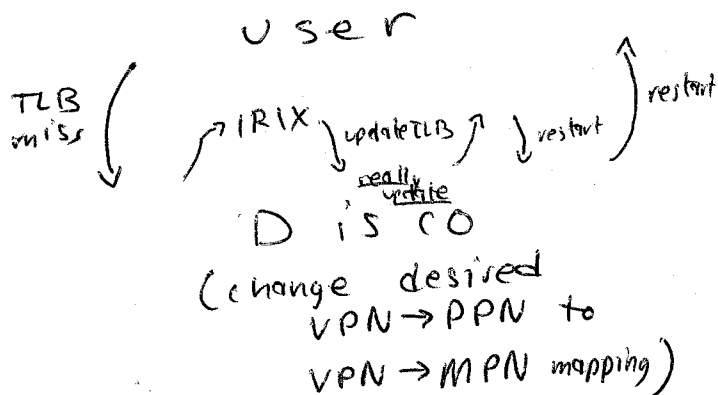
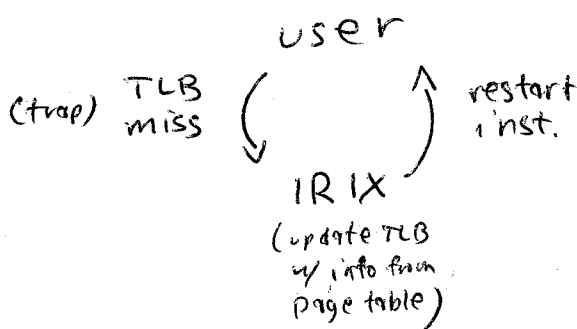
\Rightarrow to speed up: TLB

- in hardware
- fully associative cache of VA \rightarrow PA

TLB miss: "software handled" TLB

usual

Disco



Memory Difficulties

5

⇒ IRIX usually lives in unmapped physical memory?
(KSEG 0)

⇒ what does this mean?
(can't interpose via TLB)

⇒ have to relink kernel
(now KTLB faults an issue)

⇒ ASID (address space identifier)

⇒ what is it?

⇒ Hard to virtualize

⇒ Flush TLB on VCPU switch

⇒ Cost of TLB miss is high

⇒ Add 2nd-level S/W TLB

NUMA issues

6

Key: w/ $VA \rightarrow PA \rightarrow MA$,
can move "phys" pages to handle NUMA issues

Goal: cache misses serviced from local memory

Example: when to replicate a page?

\Rightarrow Read sharing (e.g. kernel code)

Example: when to migrate a page?

\Rightarrow when write activity occurring from a far away CPU (e.g. scheduler moved VM)

How to do these things?

\Rightarrow must modify TLB appropriately

migrate: invalidate

replicate: downgrade to read-only

What about heavily write-shared pages?

\Rightarrow don't move, no point

Handling I/O (SKIP)

(7)

PIO:

usual

memory map control registers
use loads/stores to interact
w/ device

disco

add special driver to OS
(network, SCSI)

use internal "monitor call" interface
to interact w/ I/O device efficiently

DMA

usual

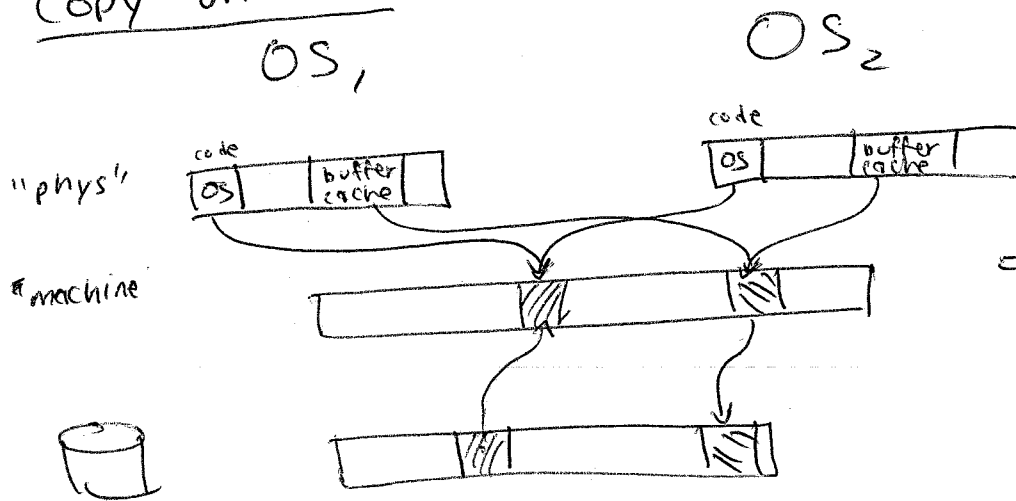
for efficient, large transfers w/o CPU

disco

must intercept, translate $PA_s \rightarrow MA_s$

Sharing across VMs

Copy-on-write

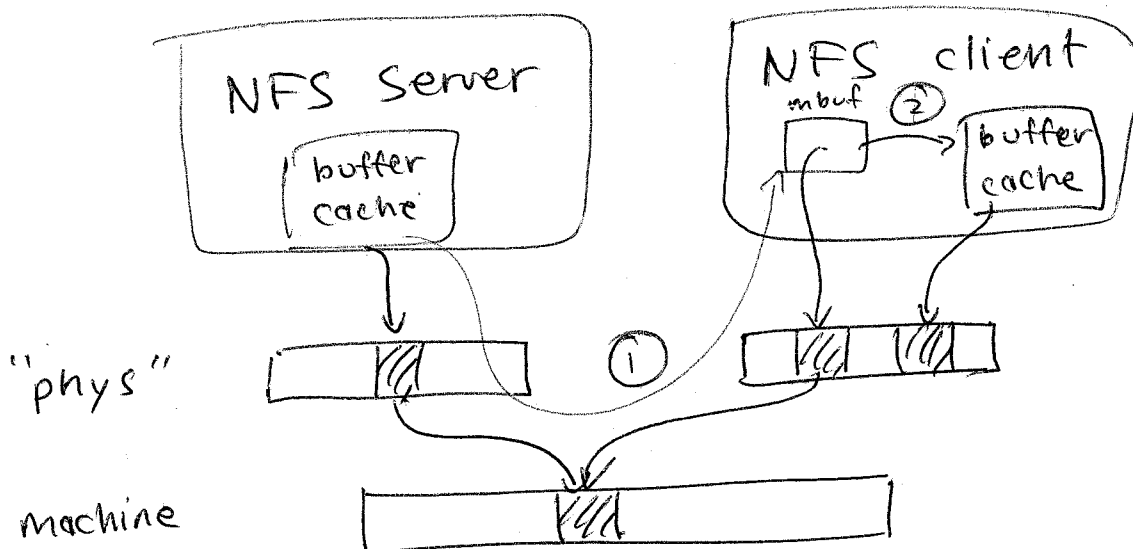


⇒ intercept DMA from block X;
if already in memory, just
remap PA → MA

[COW]

⇒ good for read-only sharing
OS text, some of buffer cache

Network interface



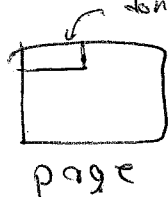
- ① send becomes remapping
- ② but client still copies incoming packet
(solve by changing client)

Breaking the Illusion

9

- ① some priv. ops just need to read/write regs
⇒ replace w/ LD/ST to mem addresses
- ② zeroed pages:
⇒ why zero @ all?
⇒ Problem: OS zeroes, monitor zeroes
⇒ Solution: OS asks monitor for zeroed page
- ③ Page on free list
⇒ OS calls monitor to tell it
- ④ CPU is idle
⇒ Disco detects low power mode
⇒ idleness
- ⑤ bcopy → remap in NFS client

- ⑥ mbuf structure change



Performance #5

10

⇒ All simulated results
(but reliable)

⇒ microbenchmarks : exec,
open, write 1.6 → 2x slower

⇒ macro : apps run well

Conclude

⇒ Power of layering , illusion

⇒ Importance of info flow

⇒ Today : VMware

key : consolidation in server farms