

## **Disco: Running Commodity Operating Systems on Scalable Multiprocessors**

**Bugnion, Devine, and Rosenblum**

**SOSP 1997**

1. What was the goal of Disco?

**Disco: Running Commodity Operating Systems on Scalable Multiprocessors**  
**Bugnion, Devine, and Rosenblum**  
**SOSP 1997**

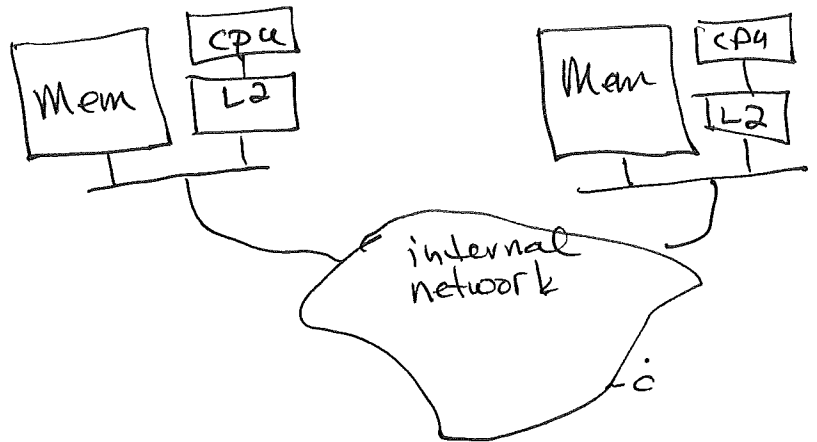
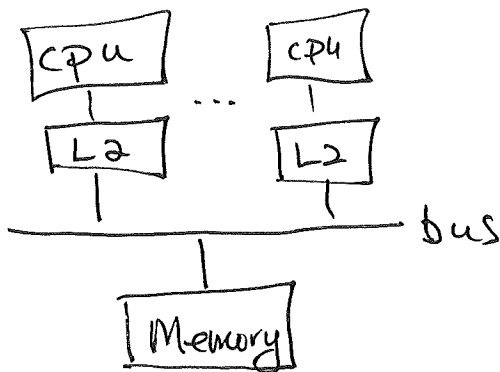
1. What was the goal of Disco?

Run commodity OS on NUMA

2. What are the pros and cons of a CC-NUMA architecture (vs SMP UMA)? What are some issues with building or modifying an OS for CC-NUMA?

2. What are the pros and cons of a CC-NUMA architecture (vs SMP UMA)? What are some issues with building or modifying an OS for CC-NUMA?

## cache-coherent Non-Uniform Memory Architecture



+ same performance  
to all memory  
from everywhere

- crazy perf  
+ scalable hw

- hard to make hw  
scalable - hit limits  
of bus

(correctness the same)

---

o) Write from scratch

Difficulties:

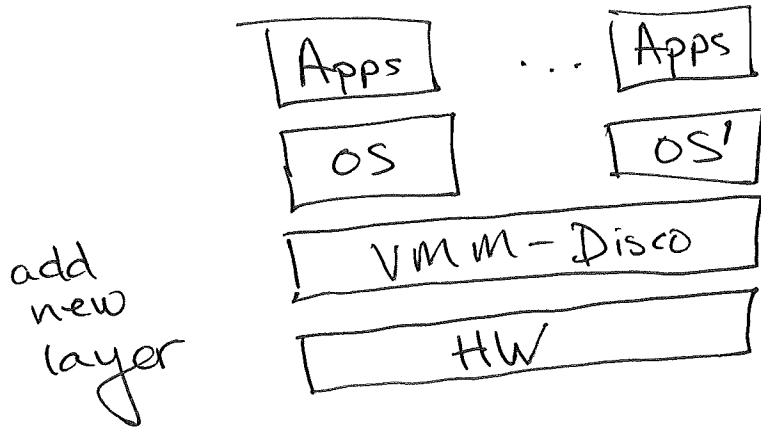
1) Take existing OS + port: huge, ugly, difficult to understand

- need to allocate OS data structures  
across NUMA (hard enough for  
SMPs)  
locks

⇒ VMM!

3. What are the advantages of Virtual Machines?

### 3. What are the advantages of Virtual Machines?



- hide tough issues from OS
- portability layer
- smaller, easier to understand + trust
- run different OSes concurrently (almost unmodified)

4. At a high level, what are some of the challenges of Virtual Machine?

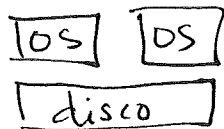
4. At a high level, what are some of the challenges of Virtual Machine?

### 1) Overhead - cost of virtualizing

a) Time: 
$$\frac{\frac{\text{app}}{\text{OS}}}{\frac{\text{disco}}{\text{hw}}}$$

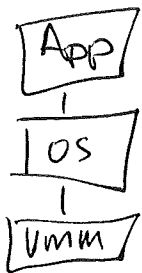
vmm acts as emulator  
- most instr can just run,  
but priv. instr. + TLB inst  
must be trapped +  
emulated

b) Space:



- multiple copies wastes memory  
- OS code + file cache

### 2) Resource management



- lose information about what is being used

a) cpu - idle thread?

b) memory - free list

### 3) Sharing problems

- most OS require exclusive access to disk



5. How does Disco virtualize the MIPS R10000 CPU? What happens on a system call w/o and w/ Disco? Why are 3 modes for user, supervisor, and kernel key?

a  
c  
5. How does Disco virtualize the MIPS R10000 CPU? What happens on a system call w/o and w/ Disco? Why are 3 modes for user, supervisor, and kernel key? b

-most: direct execution

-track each VM that must be run,  
set real registers to VCPU regs, jmp to VPC (similar to OS + process table entries)

-hard? priv. instr.

MIPS detail:

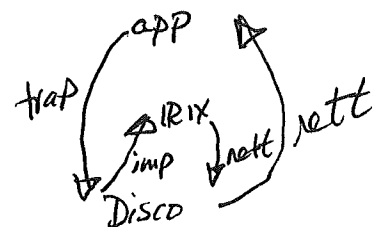
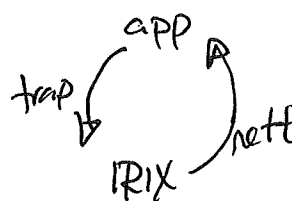
user  
supervisor  
kernel

-access all segments of mem, but not priv. instr, or phys mem

usual  
app  
-  
IRIX

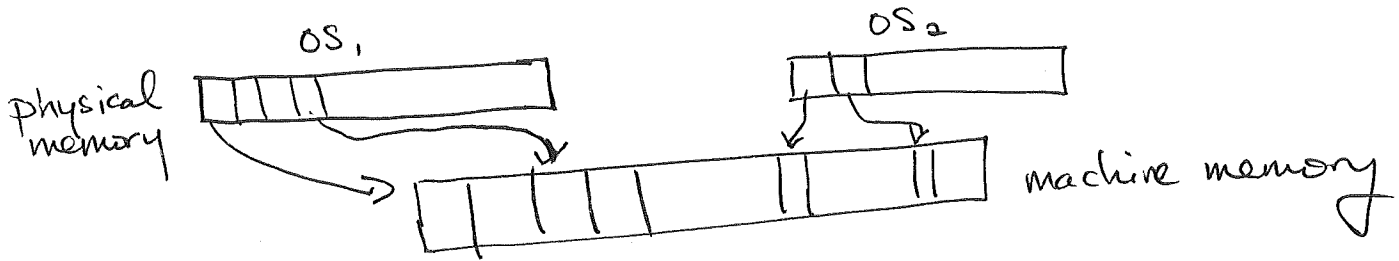
Disco  
app  
OS-IRIX  
DISCO

OS can still access all of memory + manage apps



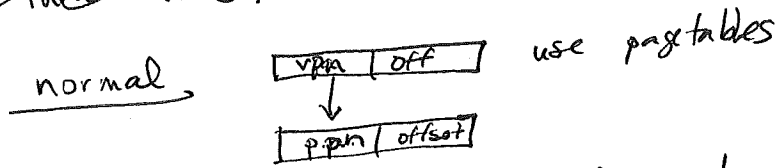
6. How does Disco virtualize memory? What will be held in the TLB? What happens on a TLB miss with w/o and w/ Disco? What data structure does Disco add?

6. How does Disco virtualize memory? What will be held in the TLB? What happens on a TLB miss with w/o and w/ Disco? What data structure does Disco add?

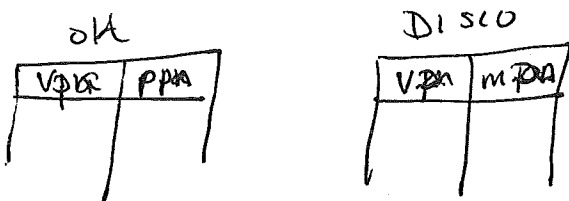


V.A. → P.A. → M.A.

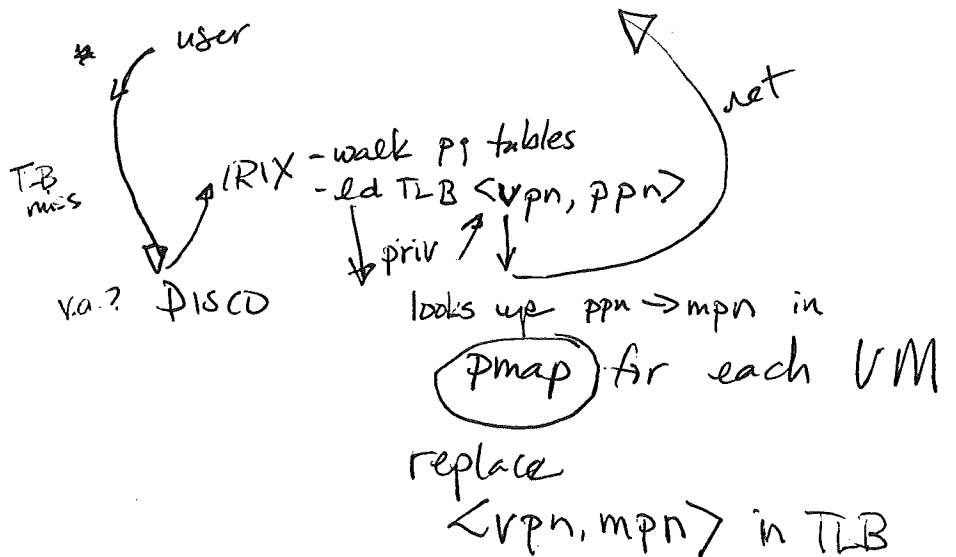
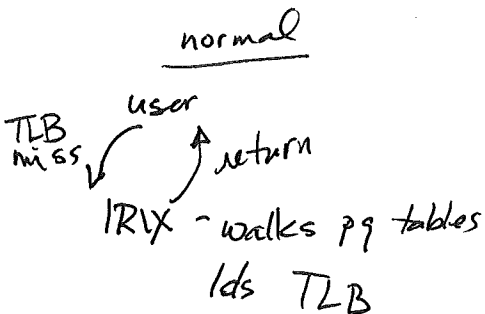
~~Do then TLB:~~



could trap on every ld + st,



TLB miss?



7. What complexity was caused by IRIX living in kseg0? What was Disco's solution?

7. What complexity was caused by IRIX living in kseg0? What was Disco's solution?

KSEG0 - unmapped physical memory

→ can't interpose via TLB

- all OS memory references ~~would be~~ ...

→ Relink IRIX to use different segment

8. Why are TLB misses more significant with Disco? What is Disco's solution?

8. Why are TLB misses more significant with Disco? What is Disco's solution?

1) more costly

2) more of them

a) OS is using TLB too

b) must flush TLB between VCPU switches; Why?

TLB

ASID	vpn	mpn
------	-----	-----

difficult/costly to virtualize ASID across all VMs, so instead flush

⇒ Add 2nd-level TLB in SW

to do quick replacements - no need to walk pt in IPX  
(modify prev. diagram)



9. What are Disco's goals that are specific to NUMA? When should a page be replicated? migrated? How does Disco perform replication and migration? What should happen if a page is heavily write-shared?

9. What are Disco's goals that are specific to NUMA? When should a page be replicated? <sup>rep</sup>migrated? How does Disco perform <sup>rep</sup>replication <sup>rep</sup>and <sup>rep</sup>migration? What should happen if a page is heavily write-shared?

- Handle performance of NUMA (not correctness)
    - cache misses from CPU handled in local memory
- Key: Use  $VA \rightarrow PA \rightarrow MA$  to handle issues

→ - Replicate?

- Read sharing ~~(OS text)~~

→ - Migrate?

- Activity on new CPU (e.g. <sup>Disco</sup>scheduler moved VM) w/ affinity, tries not to, but...

How?

- modify old TLB appropriately
- migration: invalidate! old mpu entries
- replication: downgrade to read-only

|  
If heavily-write-shared?

Don't migrate or replicate

10. Why are large memory footprints a concern for Disco? Why does sharing occur across VMs? Why is copy-on-write useful?

10. Why are large memory footprints a concern for Disco? Why does sharing occur across VMs? Why is copy-on-write useful?

- Multiple OSes now running; encourage use of NFS server for shared file access  
text-code  
└ data
- OSes share same code; NFS shares file data across client + server

Figure 3.

- intercept DMA from block X
- if already in memory, just record PA → MA
- record copy-on-write to track shared data efficiently

Figure 4.

- ① send becomes additional mapping (device) emulate
- ② copy becomes add. mapping

11. Running a completely unmodified commodity OS on Disco is tricky. What changes did Disco make to IRIX to improve performance?

11. Running a completely unmodified commodity OS on Disco is tricky. What changes did Disco make to IRIX to improve performance?

① Some priv. ops are very simple (reading regs)  
→ replace w/ ld/st to mem. addresses

② zero'd pages:

- ~~both~~ OS does when allocates new page to process (must for privacy)
- Disco must across VMs
- avoid double work

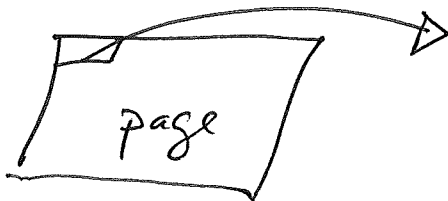
③ Pages on free list - Disco should know about

④ CPU is idle

- Disco detects low power

⑤ bcopy → remap in NFS client

⑥ mbuf structure



12. As shown in Figure 5, how much time overhead does DISCO impose for a uniprocessor workload? Why does some of the original kernel time decrease?

12. As shown in Figure 5, how much time overhead does DISCO impose for a uniprocessor workload? Why does some of the original kernel time decrease?

- more for compute-bound w/ TLB misses
- pmake uses many system services
- Zeroing pages
- 2nd-level TLB



13. What does Figure 6 show? Does Disco do a decent job sharing buffer cache space across VMs? Of sharing IRIX text? IRIX data?

13. What does Figure 6 show? Does Disco do a decent job sharing buffer cache space across VMs? Of sharing IRIX text? IRIX data?

yes

yes

no

- pmake workload
- virtual footprint vs. machine memory

14. What does Figure 7 show? Where can you find an evaluation of Disco's replication and migration policies?

14. What does Figure 7 show? Where can you find an evaluation of Disco's replication and migration policies?

8-processor cc-NUMA

1P1X · Lots of time in kernel <sup>memlock scalability limit</sup>

1VM · just add DISCO overhead

2VM → 8VM actually improves  
because DISCO does not have bad lock

turned off migration + replication

---

Fig. 8.

- much less time accessing remote mem

## 15. Conclusions?

## 15. Conclusions?

### Strengths

- + Power of layering / indirection
- + small monitor hides tough NUMA + parallel issues
- + Really works - useful in many settings (not just cc-NUMA)

### Problems:

- True virtualization tricky
- Handling into loss than stack