

Lecture #11

CS695

Memory Resource Management

- primitives
- mgmt. policies / decisions.
 - content-based

ballooning, sharing, migration,

check pointing,

record & replay

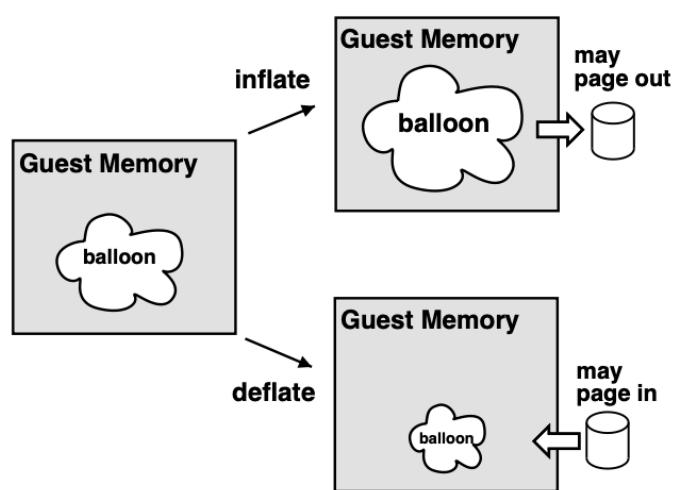
~ symbiotic: pressure by NMM, eviction by guest OS

~ pinning

↳ to avoid reclaimation

~ all ballooned pages
in guest valid=1 } pte
Present=0 } flag

bits



Q. Is this experiment reproducible?

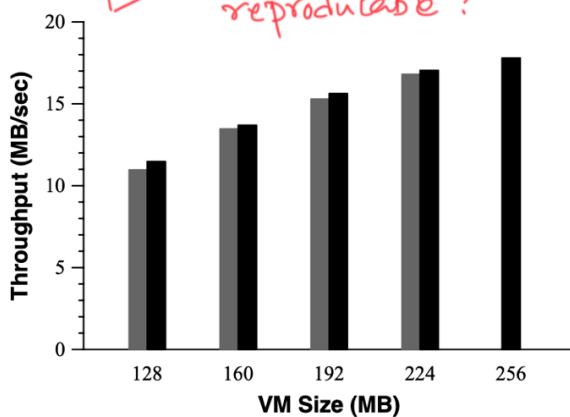


Figure 2: **Balloon Performance.** Throughput of single Linux VM running dbench with 40 clients. The black bars plot the performance when the VM is configured with main memory sizes ranging from 128 MB to 256 MB. The gray bars plot the performance of the same VM configured with 256 MB, ballooned down to the specified size.

What was workload pattern?

of each of the

& 0 clients.

or removed from a VM in order to rapidly adjust its

Sharing (content based page sharing)

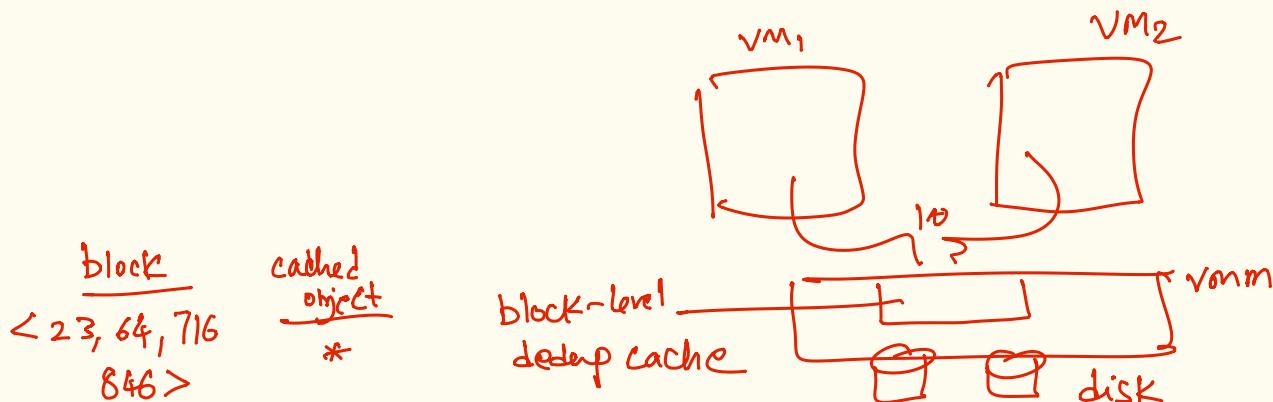
~ same memory regions / pages across VMs is a common possibility

Solⁿ: ~~try~~ deduplicate!

~ processes: code, data, heap, stack.

~ VMM-view: ppa \rightarrow mpa (no semantics of page use)

Solⁿ: ~ intercept disk I/O calls (device virtualization)



memory sharing of VMs.

~ identify similar pages by content

~ hypervisor has access to all P2M mappings of all VMs

~ scan pages periodically / background...

~ for each scanned page

~ hash the page

~ lookup if hash exists

~ if yes

- ① recompute hash of hint page & compare
 - ② byte-by-byte compare of pages. { $P_1 \rightarrow m$ } { $P_2 \rightarrow m$ }
- } share a m/c page
- mark as cold
 - ref count

else, add hash to a hint set

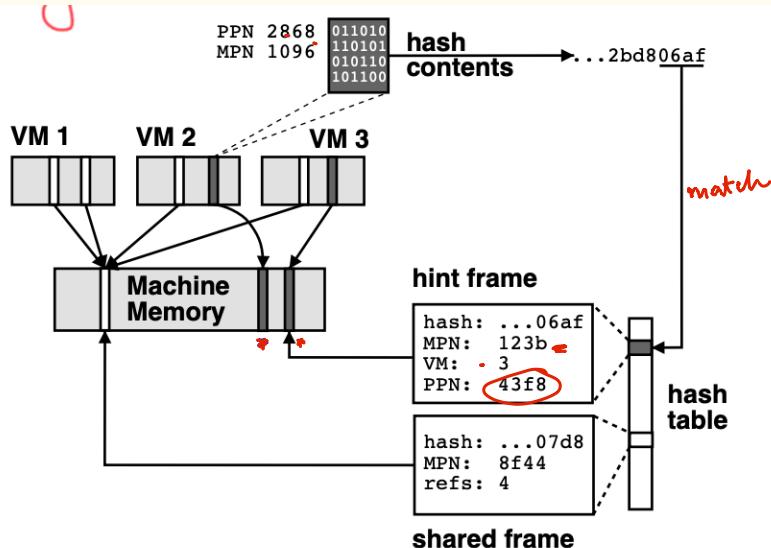


Figure 3: Content-Based Page Sharing. ESX Server scans for sharing opportunities, hashing the contents of candidate PPN 0x2868 in VM 2. The hash is used to index into a table containing other scanned pages, where a match is found with a hint frame associated with PPN 0x43f8 in VM 3. If a full comparison confirms the pages are identical, the PPN-to-MPN mapping for PPN 0x2868 in VM 2 is changed from MPN 0x1096 to MPN 0x123b, both PPNs are marked COW, and the redundant MPN is reclaimed.

Xinru
- KSM — Kernel same page merge
- madvise

questions CBPS operation

- + - when to scan?
- how much to scan?
- where to scan?

time to discover sharing potential?

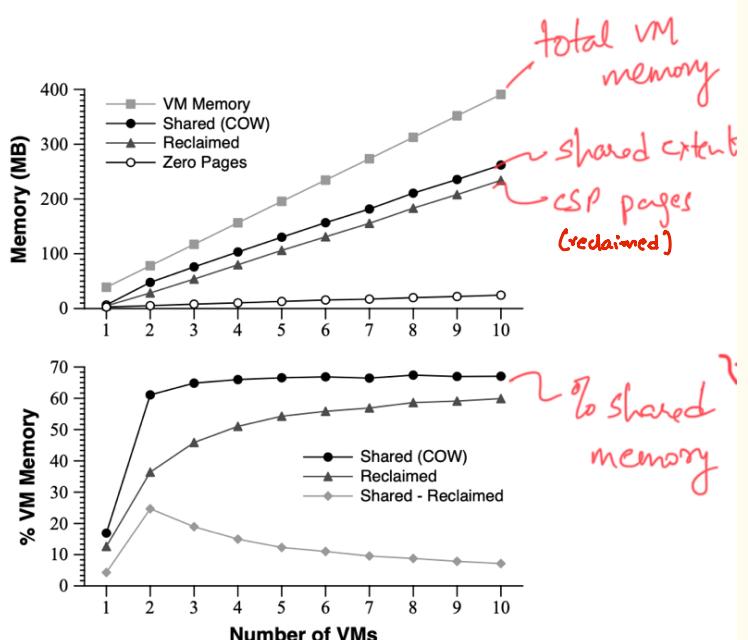


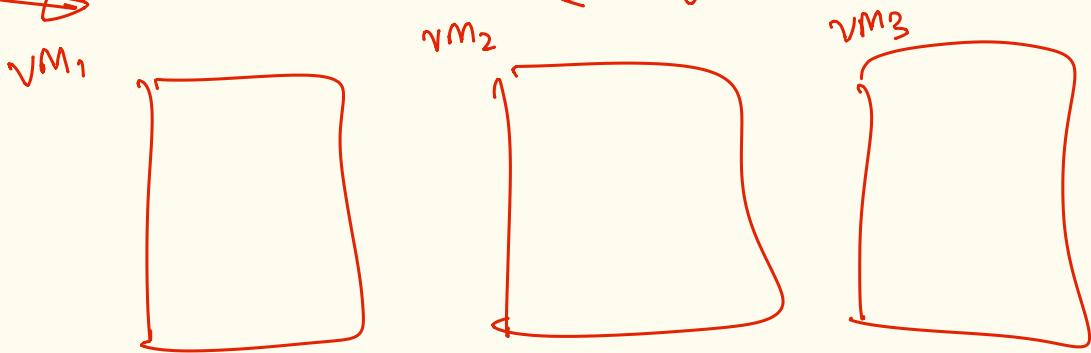
Figure 5: Real-World Page Sharing. Sharing metrics from production deployments of ESX Server. (a) Ten Windows NT VMs serving users at a Fortune 50 company, running a variety of database (Oracle, SQL Server), web (IIS, Websphere), development (Java, VB), and other applications. (b) Nine Linux VMs serving a large user community for a nonprofit organization, executing a mix of web (Apache), mail (Major-domo, Postfix, POP/IMAP, MailArmor), and other servers. (c) Five Linux VMs providing web proxy (Squid), mail (Postfix, RAV), and remote access (ssh) services to VMware employees.

Figure 4: Page Sharing Performance. Sharing metrics for a series of experiments consisting of identical Linux VMs running SPEC95 benchmarks. The top graph indicates the absolute amounts of memory shared and saved increase smoothly with the number of concurrent VMs. The bottom graph plots these metrics as a percentage of aggregate VM memory. For large numbers of VMs, sharing approaches 67% and nearly 60% of all VM memory is reclaimed.

order of usage of primitives for memory mgmt.

- (i) sharing
- (ii) ballooning
- (iii) demand paging/ swapping.

how to decide memory allocation extents?



balance performance & memory efficiency/utility.

Scheme : WSS — working set size
↳ active memory region / area

Share-based allocation
proportionate

$$\text{effective } \rightarrow f = \frac{s}{p \cdot (f + k(1-f))}$$

prioring / fraction
for allocation

$$f_1 : 0.3$$

$$f_2 : 0.2$$

$$f_3 : 0.1$$

$$S = \# \text{shares}/\text{currency of a VM}$$

$$0 < T \leq 1$$

$$P = \# \text{pages}$$

$$f = \text{fraction of active pages}$$

$$k = \frac{1}{1-T}$$

K = in-active tax.

$$T = \text{tax rate } g = S/p$$

$$T = 0 \Rightarrow k = 1$$
$$= 1 \Rightarrow k = \infty \quad g = 0$$

#

sampling + ~~minor~~
~~page~~
 fault

randomly mark n pages

invalid

$P=0$

$v=1$

} page fault

count 'k' such page faults

$$f \approx \frac{k}{n}$$

periodic
 loop