

CITRIX[®]

Xen and XenServer Storage Performance

Low Latency Virtualisation Challenges

Dr Felipe Franciosi XenServer Engineering Performance Team

e-mail: felipe.franciosi@citrix.comfreenode: felipef #xen-apitwitter: @franciozzy

Agenda

- Where do Xen and XenServer stand?
 - When is the virtualisation overhead most noticeable?
- Current implementation:
 - blkfront, blkback, blktap2+tapdisk, blktap3, qemu-qdisk
- Measurements over different back ends
 - Throughput and latency analysis
 - Latency breakdown: where are we losing time when virtualising?
- Proposals for improving



When is the virtualisation overhead noticeable?



- What kind of throughput can I get from dom0 to my device?
 - Using 1 MiB reads, this host reports 118 MB/s from dom0

File View Pool Server VM Store	nge Templates Tools Window Help w Server 🏪 New Pool 🛅 New Storage 🛅 New VM 🥹 Shut Down 🛞 Reboot 🕕 Suspend 🛛 ✔ No System Alerts	
Views: Server View	togged in as: Local root account	
Search 🔎	Search General Memory Storage Networking NICs Console Performance Users Logs	
XenCenter	dt01 server console	f=/dev/null bs=1M count=500 iflag=direct , 4.44602 seconds, 118 MB/s

- What kind of throughput can I get from domU to my device?
 - Using 1 MiB reads, this host reports 117 MB/s from a VM

5

File Were Pool Server WM Storage Templete Tools Window Help Image: File Image: File <th>🛿 XenCenter</th> <th></th> <th></th> <th></th>	🛿 XenCenter			
Construct → Construct → Add New Storer → New Storage → New VM → State → Storage → New Stora	File View Pool Server VM St	orage Templates Tools Window Help		
Send Ctrl+Alt+Del (Ctrl+Alt+Insert) Scale 🖉 Undock (Alt+Shift+U) Fullscreen (Ctrl+Enter)	File View Pool Server VM Str Search Search Search Search Search Search Search Search Search Search Search Search <t< th=""><th>Orage Templates Tools Window Heip New Server Image: Server New Yorage New Storage New VM Image: Server <t< th=""><th>Shut Down ♀ Reboot ① Suspend ✓ No System Alerts Logged in as: Local root account ance Snapshots Logs ✓ Eject Looking for guest console /dev/null bs=1M count=500 if lag=direct 5319 s, 117 MB/s Eezy32: ~ # dd if =/dev/xvdb ecords in ecords out 300 bytes (524 MB) copied, eezy32: ~ # ■</th><th>IMPERCEPTIBLE virtualisation overhead of=/dev/null bs=1M count=500 iflag=direct 4.46319 s, 117 MB/s</th></t<></th></t<>	Orage Templates Tools Window Heip New Server Image: Server New Yorage New Storage New VM Image: Server Image: Server <t< th=""><th>Shut Down ♀ Reboot ① Suspend ✓ No System Alerts Logged in as: Local root account ance Snapshots Logs ✓ Eject Looking for guest console /dev/null bs=1M count=500 if lag=direct 5319 s, 117 MB/s Eezy32: ~ # dd if =/dev/xvdb ecords in ecords out 300 bytes (524 MB) copied, eezy32: ~ # ■</th><th>IMPERCEPTIBLE virtualisation overhead of=/dev/null bs=1M count=500 iflag=direct 4.46319 s, 117 MB/s</th></t<>	Shut Down ♀ Reboot ① Suspend ✓ No System Alerts Logged in as: Local root account ance Snapshots Logs ✓ Eject Looking for guest console /dev/null bs=1M count=500 if lag=direct 5319 s, 117 MB/s Eezy32: ~ # dd if =/dev/xvdb ecords in ecords out 300 bytes (524 MB) copied, eezy32: ~ # ■	IMPERCEPTIBLE virtualisation overhead of=/dev/null bs=1M count=500 iflag=direct 4.46319 s, 117 MB/s
	© 2013 Citrix		***	CİTRI

• That's not always the case...

6

Same test on different hardware (from dom0)



• That's not always the case...

7

Same test on different hardware (from domU)



Current Implementation

How we virtualise storage with Xen



Current Implementation (bare metal)

• How does that compare to storage performance again?

- There are different ways a user application can do storage I/O
 - We will use simple read() and write() libc wrappers as examples



Current Implementation (Xen)

- The "Upstream Xen" use case
 - The virtual device in the guest is implemented by blkfront
 - Blkfront connects to blkback, which handles the I/O in dom0



Current Implementation (Xen)

- The XenServer 6.2.0 use case
 - XenServer provides thin provisioning, snapshot, clones, etc. hello VHD
 - This is easily implemented in user space. hello TAPDISK



Current Implementation (Xen)

- The blktap3 and qemu-qdisk use case
 - Have the entire back end in user space



Throughput and latency analysis



- Same host, different RAID0 logical volumes on a PERC H700
 - All have 64 KiB stripes, adaptive read-ahead and write-back cache enabled















• There is another way to look at the data:



1

= Latency (time/data)







- Initial tests showed there is a "warm up" time
- Simply inserting printk()s affect the hot path
- Used a trace buffer instead
 - In the kernel, trace_printk()
 - In user space, hacked up buffer and a signal handler to dump its contents
- Run 100 read requests
 - One immediately after the other (requests were sequential)
 - Used IO Depth = 1 (only one in flight request at a time)
- Sorted the times, removed the 10 (10%) fastest and slowest runs
- Repeated the experiment 10 times and averaged the results





• The penalty of copying can be worth taking depending on other factors:

(TLB flushes)

(contention on grant tables)

- Number of dom0 vCPUs
- Number of concurrent VMs performing IO
- Ideally, blkfront should support both data paths
- Administrators can profile their workloads and decide what to provide

Proposals for Improving

What else can we do to minimise the overhead?

• How can we reduce the processing required to virtualise I/O ?

27 © 2013 Citrix

- Persistent Grants
 - Issue: grant mapping (and unmapping) is expensive
 - Concept: back end keeps grants, front end copies I/O data to those pages
 - Status: currently implemented in 3.11.0 and already supported in gemu-gdisk

Indirect I/O

- Issue: blkif protocol limit requests to 11 segs of 4 KiB per I/O ring
- Concept: use segments as indirect mapping of other segments
- Status: currently implemented in 3.11.0

- Avoid TLB flushes altogether (Malcolm Crossley)
 - Issue:
 - When unmapping, we need to flush the TLB on all cores
 - But in the storage data path, the back end doesn't really access the data
 - Concept: check whether granted pages have been accessed
 - Status: early prototypes already developed
- Only grant map pages on the fault handler (David Vrabel)

• Issue:

- Mapping/unmapping is expensive
- But in the storage data path, the back end doesn't really access the data
- Concept: Have a fault handler for the mapping, triggered only when needed
- Status: idea yet being conceived

29 © 2013 Citrix

CİTRİX®

Split Rings

• Issue:

- blkif protocol supports 32 slots per ring for requests or responses
- this limits the total amount of outstanding requests (unless multi-page rings are used)
- this makes inefficient use of CPU caching (both domains writing to the same page)
- Concept: use one ring for requests and another for responses
- Status: early prototypes already developed
- Multi-queue Approach
 - Issue: Linux's block layer does not scale well across NUMA nodes
 - Concept: allow device drivers to register a request queue per core
 - Status: scheduled for Kernel 3.12

e-mail: felipe.franciosi@citrix.comfreenode: felipef #xen-apitwitter: @franciozzy