

# Scale 14x: The Latest From the Xen Project

Lars Kurth

Cat Herder, Xen Project

Chairman, Xen Project Advisory Board



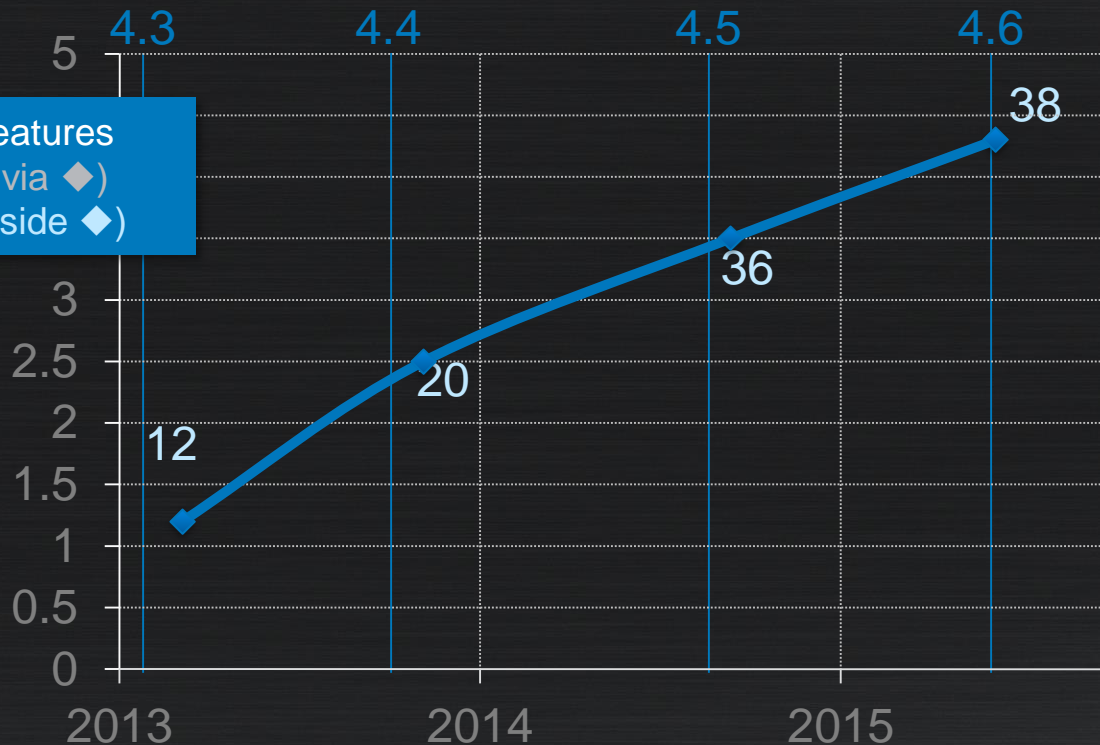
lars\_kurth



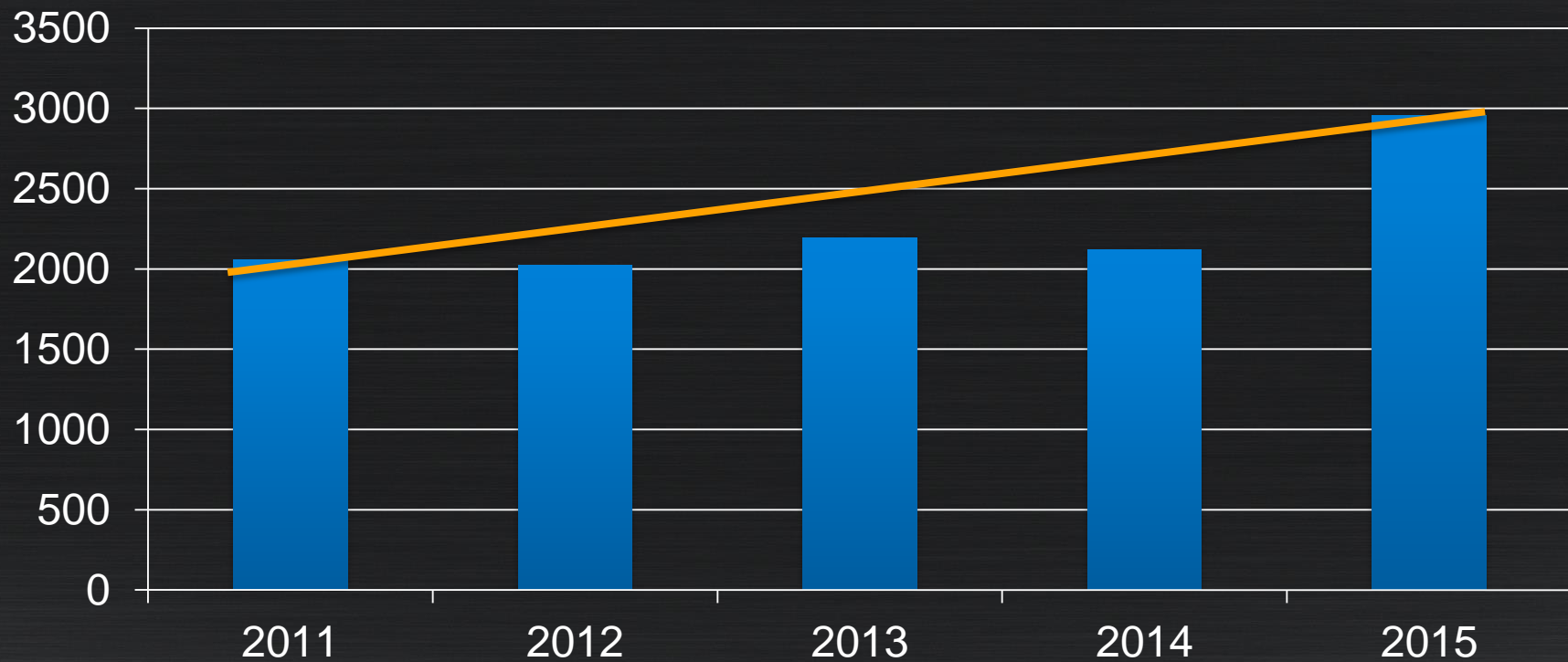
# Xen 4.x Hypervisor Release History

## Xen 4.x Number of New Major Features

- Developed per month (y axis via ◆)
- Absolute number (number beside ◆)



# Hypervisor Git Commits





# Rate of Innovation is Accelerating

While quality and security requirements are increasing simultaneously

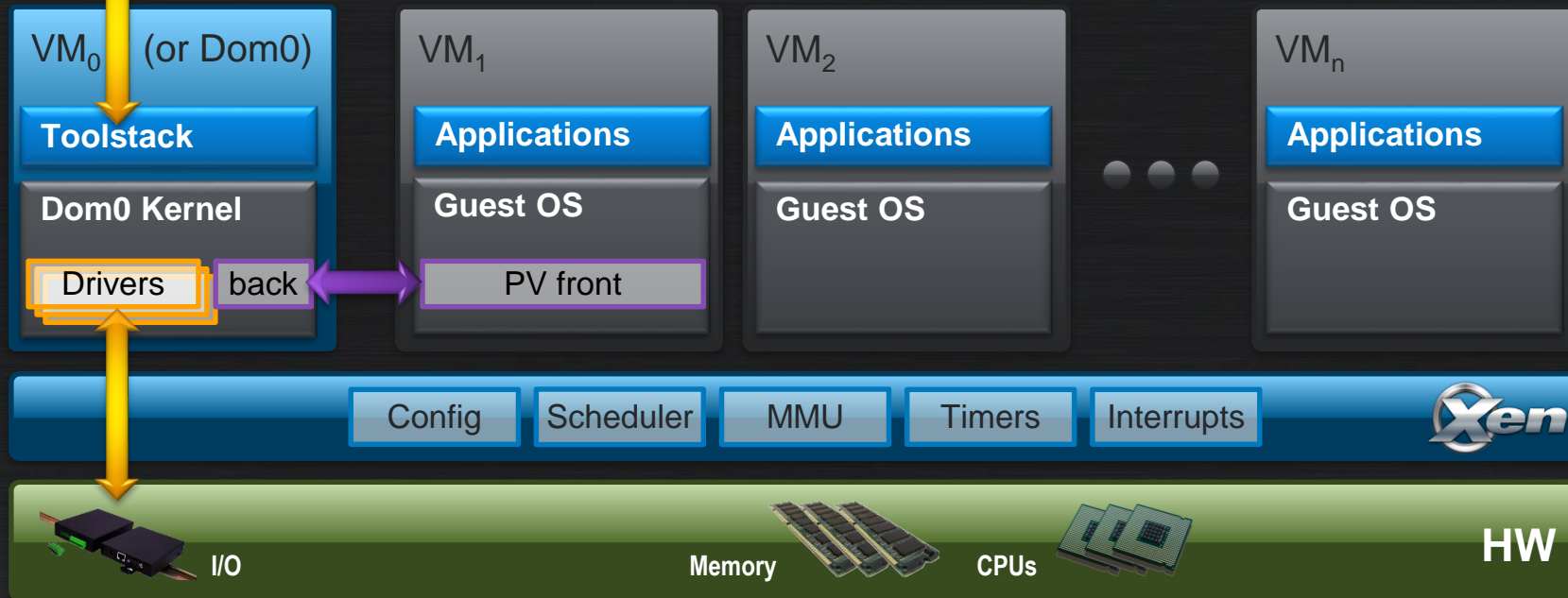


# A type-1 Hypervisor

with a twist



Console





# Architectural Advantages

**Density:** It's thin

Excellent for supporting many very small workloads (e.g. unikernels)

**Scalability:** It can support huge numbers of VMs

Terrific for highly dense workloads (e.g. unikernels, disaggregation, ...)

**Security:** Host OS isolated within a VM

This makes it harder to attack the Host OS

**Scheduling:** Can use dedicated scheduler

Enables specialized workload profiles (mix and match schedulers on one host)

**Paravirtualization:** Simplified interface

Easy to implement a unikernel base

Enables fast boot times necessary for unikernels



# The Interface Between Xen and Unikernels

Virtualisation Modes  
The Future of Virtualization Modes  
Implications for Unikernel bases



# Evolution of Virtualization Modes (x86)

Shortcut	Xen	Mode	With
HVM / Fully Virtualized		HVM	
HVM + PV drivers	3.0	HVM	PV Drivers
PVHVM	4.0	HVM	PVHVM Drivers
PVH	4.4/4.5	PV	pvh=1
PV		PV	

## Unikernel Bases:

Primarily depend on PV  
E.g. rumprun and Mini-OS  
Will work on Xen based  
clouds and hosting services

	Poor Performance
	Scope for Improvement
	Optimal Performance

PV = Paravirtualized

VS = Software Virtualized

VH = Hardware Virtualized

Shortcut	Mode	With				
			Disk and Network	Interrupts & Timers	Emulated Motherboard, Legacy Boot	Privileged Instructions, Page Tables
HVM / Fully Virtualized	HVM		VS	VS	VS	VH
HVM + PV drivers	HVM	PV Drivers	PV	VS	VS	VH
PVHVM	HVM	PVHVM Drivers	PV	PV	VS	VH
PVH	PV	pvh=1	PV	PV	PV	VH
PV	PV		PV	PV	PV	P

# Making PVH better

The motivation behind PVH

- HVM (like) Dom0: performance & Dom0 modification
- PVH as fast or faster than HVM
- PVH runs a PV guest within a HVM container (essentially a mix of PV & HVM)

**BUT:** PVH inherits all the PV limitations, e.g.

- Paging restrictions, lack of access to emulated devices (if needed), ...
- Concept designed prior to “additional quality and security requirements”

**Solution:** HVMLite to eventually replace PVH

- A lot simpler to implement: less code to maintain and thus to keep secure
- Behaves exactly like PVH (but internal implementation different)
- HVM without QEMU

Privileged Instructions,  
Page Tables  
Emulated Motherboard  
Legacy BIOS  
Interrupts & Timers  
Disk and Network

**Currently:**

builder = "hvm"

device\_model\_version="none"

There will need to be some  
sort of migration strategy from  
pvh=1

Shortcut	Mode	With				
HVM / Fully Virtualized	HVM				S	VH
HVM + PV drivers	HVM	PV Drivers			S	VH
PVHVM	HVM	PVHVM Drivers	PV	PV	VS	VH
<b>PVH / HVMLite</b>	HVM	pvh=1 / ...	PV	PV	PV	VH
PV	PV		PV	PV	PV	P

# Status

## Xen 4.7 (June 2016)

- HVMlite DomU support in xen.git
- Agree on config file changes and naming (PVH or HVMlite)

## Ongoing

- HVMlite Dom0 prototype for FreeBSD
- No Linux implementation yet
- Some clean-up required
- Interfaces not yet declared stable (but almost there)
- Benchmarks already very impressive



# Call to Action

PVH / HVMLite not currently used as unikernel base

- Unikernel developers make sure the architecture works for you  
(before APIs are declared stable)
- MiniOS / Rumprun not yet ported to HVMLite  
(some different approaches needed compared to pure PV)

Opportunity: Avoid Duplication

- There was a bit of duplication of unikernel bases in the early days of unikernel development (MiniOS clones)



# Performance

Focus on significant performance and scalability improvements since Xen 4.5+

# Performance Improvements

## Examples

HPET: Better and faster resolution values

Parallel memory scrubbing on boot (large machines)

Lower interrupt latency for PCI passthrough (machines > 2 sockets)

Soft affinity for non-NUMA machines

Multiple IO-REQ services for guests

(remove bottlenecks for HVM guests by allowing multiple QEMU back-ends)

SandyBridge: VT-d posted interrupts for HVM

(I/O intensive workloads)

Grant table scalability by using finer grained locks

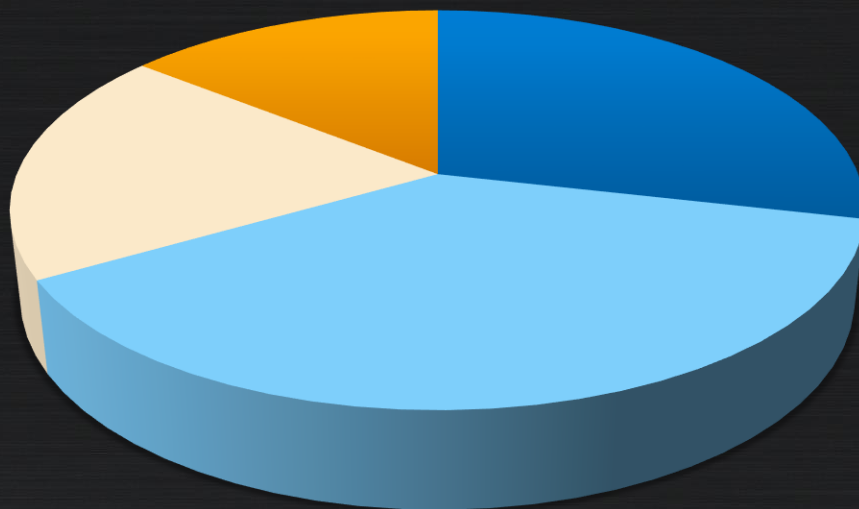
Ticket locks for improved fairness and scalability

...

# Benchmarks: Xen & KVM

(Kernel 4.2 with Xen 4.5 & QEMU 2.3)

Xen 4.6 should be even better



- Xen clearly wins
- Xen marginally wins
- KVM marginally wins
- KVM clearly wins

Source: <http://www.phoronix.com/scan.php?page=article&item=ubuntu-1510-virt>



# Schedulers

Overview  
Possibilities

## Resources:

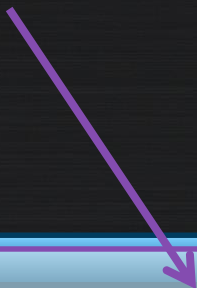
Docs: [bit.do/xen-schedulers](https://bit.do/xen-schedulers)



# Overview: Xen Project Schedulers

The Xen Project Hypervisor supports several different schedulers with different properties.

Different schedulers can be assigned to...  
... an entire host

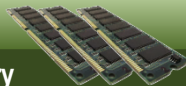


e.g. Credit2 Scheduler

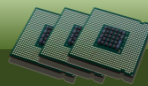


I/O

Memory



CPUs



HW

# Overview: Xen Project Schedulers

The Xen Project Hypervisor supports several different schedulers with different properties.

Different schedulers can be assigned to...

... an entire host

... a pool of physical CPU's (=CPU Pool) on a host  
(VMs need to be assigned to a pool or pinned to a CPU)



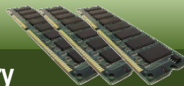
e.g. RTDS Scheduler

e.g. Credit Scheduler

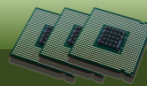


I/O

Memory

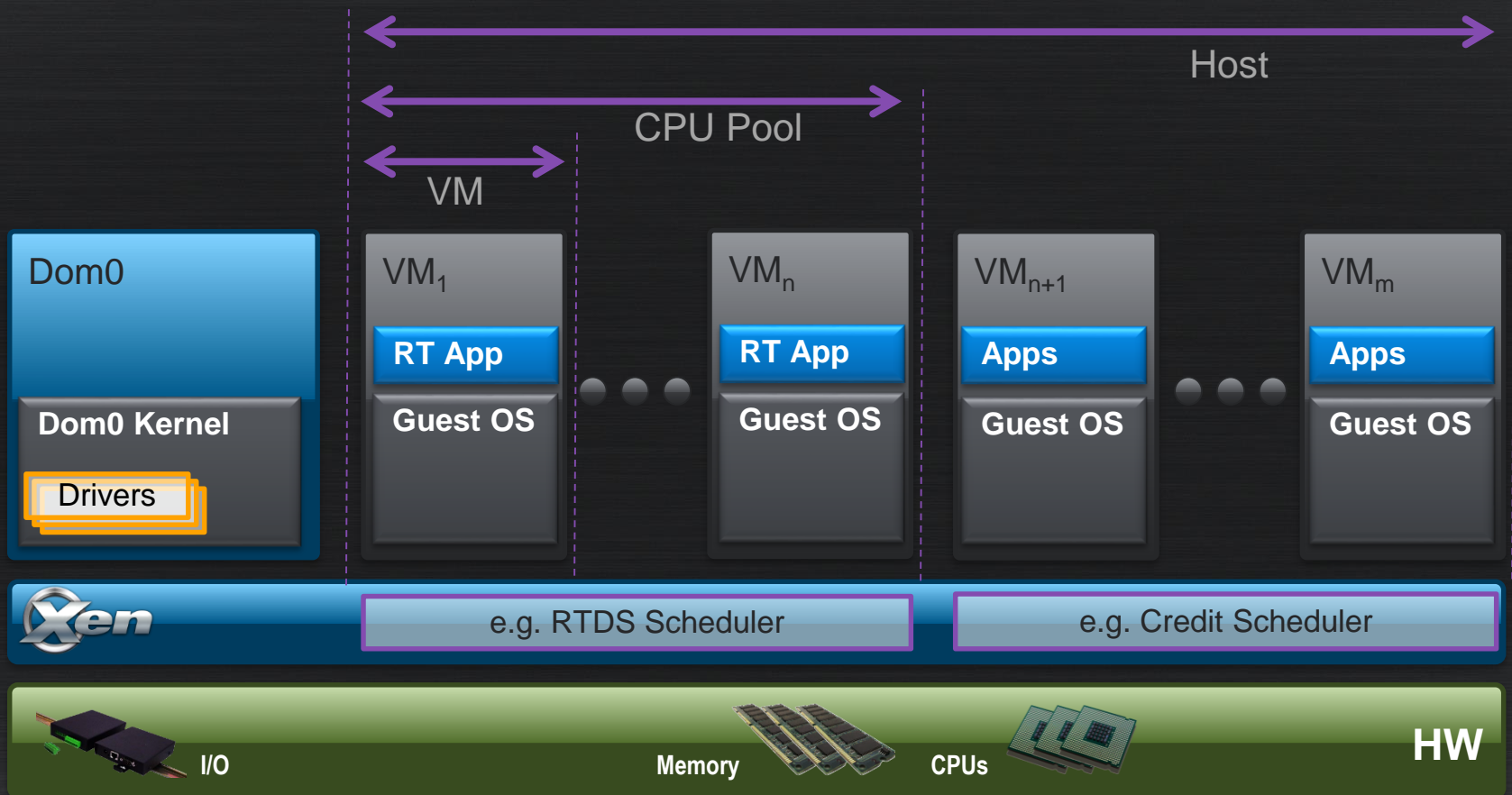


CPUs



HW

Scheduler parameters can be modified per ...



# Schedulers Overview

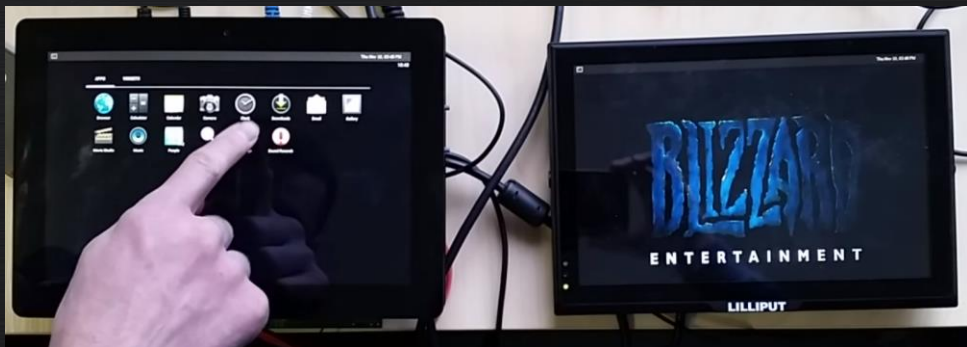
Scheduler	Use-cases	Xen 4.6	Plans for 4.7+
Credit	General Purpose	Supported <b>Default</b>	Supported <b>Default</b>
Credit 2	General Purpose Optimized for lower latency, higher VM density	Experimental	Supported
RTDS	Soft & Firm Real-time <b>Multicore</b> Embedded, Automotive, Graphics & Gaming in the Cloud, Low Latency Workloads	Hardening Optimizations Better XL support  Experimental	Adaptive granularity   Supported
ARINC 653	Hard Real-time <b>Single core</b> Avionics, Drones, Medical	Supported Compile time	No change

## Legend:

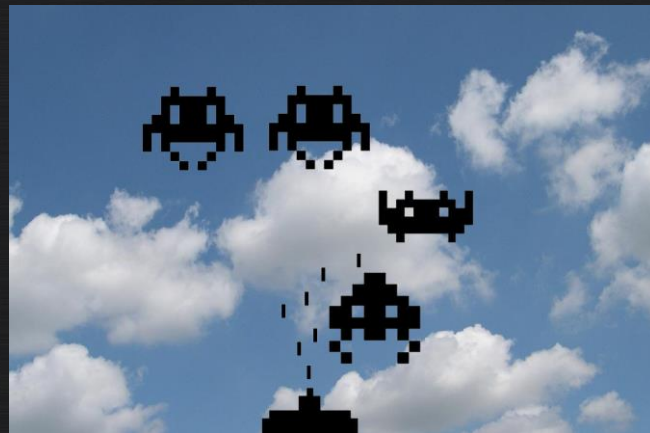
likely in 4.7

possible in 4.7

# RTDS Scheduler: Use-cases

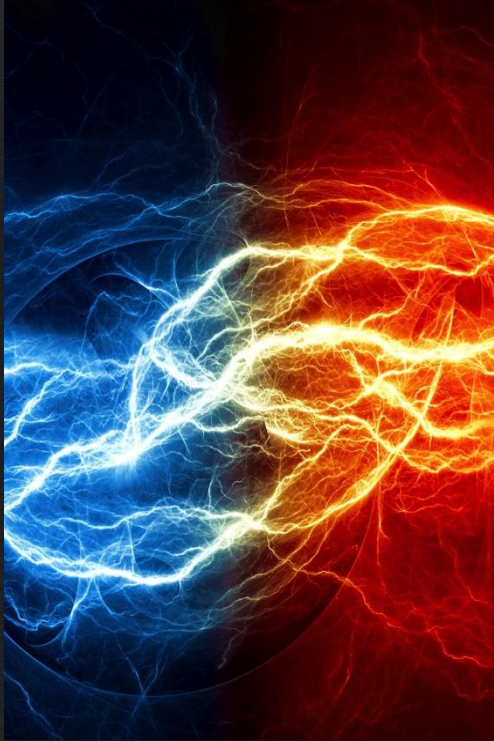


Embedded & automotive  
Latency sensitive workloads  
Guaranteed QoS



Cloud based gaming, video,  
TV delivery, ...  
Guaranteed QoS  
(Price → SLAs → QoS)





# Other major and on-going Innovations in Xen

Hardware Support

Graphics

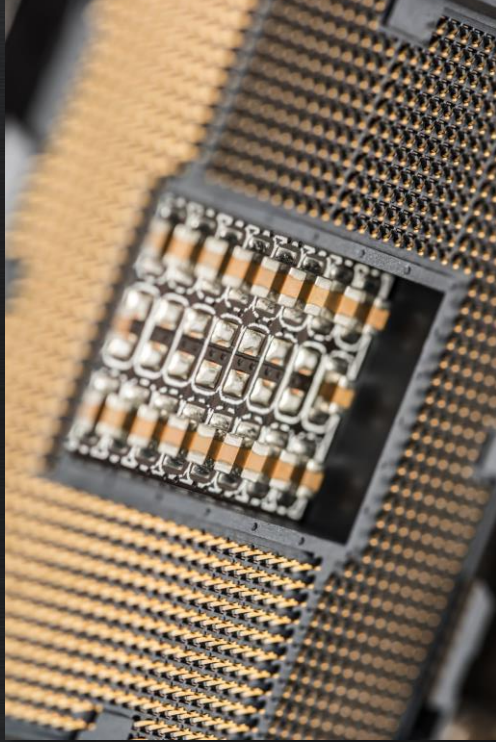
Security: VMI

Security: QEMU and Isolation

Security: xSplice

Security: Configurability

...



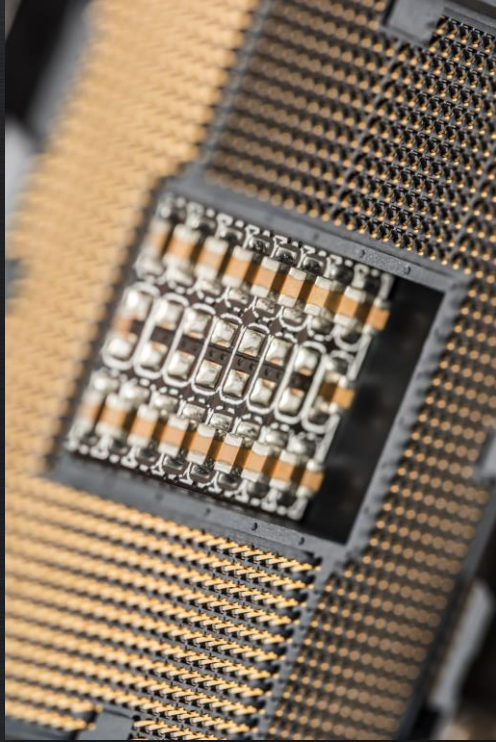
# X86 Hardware Support

Intel Platform QoS technologies  
(CMT, CAT, MBM, ...)

Virtual Performance Monitoring Unit  
vTPM v2.0

*Code/Data Prioritization*  
*Memory Protection Keys*  
*VMX TSC scaling*  
*Intel PState Drivers*  
*Posted Interrupts*

...



# ARM Hardware Support

Tracking ARM/ARM-partner server roadmap  
Hardening (more 64 bit servers in Test Lab)  
*Live Migration*



# What is Intel GVT-g (XenGT)?

## Demo Scenarios:



1. Ubuntu Dom0: 2 Windows VMs
2. VM1: World Rally Championship 3
3. VM2: Google Earth
4. VM1/VM2/Dom0 switch

Watch the demo at

<https://www.youtube.com/watch?v=V2i8HCcAnY8>

Virtual GPU per VM

Performance critical resources  
directly assigned to VM

# Intel GVT-g (XenGT) – What's next?

GVT-g support is partly out-of-tree

In use by XenClient 5.5 and XenServer Dundee

---

Most Xen patches are part of xen.git

BUT: some Linux and QEMU patches that are still in progress

Motivation: create a common code base for Xen & KVM

---

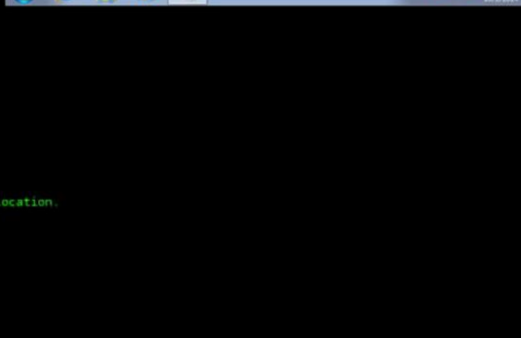
Similar approach for embedded developed by GlobalLogic  
(for ARM based architectures)



```

File Edit View Search Terminal Help
CR3=0x7b73b220 RIP=0x28685c28 ntoskrnl.exe!NtReleaseWorkerFactoryWorker
CR3=0x7b73b220 RIP=0x28685d52 ntoskrnl.exe!NtSetTimer
CR3=0x7b73b220 RIP=0x286857b1 ntoskrnl.exe!NtWaitForWorkViaWorkerFactory
CR3=0x7b73b220 RIP=0x286857b1 ntoskrnl.exe!NtWaitForWorkViaWorkerFactory
CR3=0x7b73b1e0 RIP=0x282b6435 ntoskrnl.exe!NtWaitForMultipleObjects
CR3=0x7b73b1e0 RIP=0x2824737a ntoskrnl.exe!NtClose
CR3=0x7b73b1e0 RIP=0x2824737a ntoskrnl.exe!NtClose
CR3=0x7b73b1e0 RIP=0x2825fd6d ntoskrnl.exe!NtQueryInformationThread
CR3=0x7b73b1e0 RIP=0x2824737a ntoskrnl.exe!NtClose
CR3=0x7b73b1e0 RIP=0x2824737a ntoskrnl.exe!NtClose
CR3=0x7b73b1e0 RIP=0x2824737a ntoskrnl.exe!NtClose
CR3=0x7b73b1e0 RIP=0x282b7334 ntoskrnl.exe!NtTerminateThread
CR3=0x7b73b060 RIP=0x2824737a ntoskrnl.exe!NtClose
CR3=0x7b73b060 RIP=0x28254f0a ntoskrnl.exe!NtHalpSendWaitReceivePort
CR3=0x7b73b1e0 RIP=0x2826b4db ntoskrnl.exe!NtFreeVirtualMemory
CR3=0x7b73b1e0 RIP=0x282b6435 ntoskrnl.exe!NtWaitForMultipleObjects
CR3=0x7b73b1e0 RIP=0x2824737a ntoskrnl.exe!NtClose
CR3=0x7b73b1e0 RIP=0x2824737a ntoskrnl.exe!NtClose
CR3=0x7b73b1e0 RIP=0x2825fd6d ntoskrnl.exe!NtQueryInformationThread
CR3=0x7b73b1e0 RIP=0x282b7334 ntoskrnl.exe!NtTerminateThread
CR3=0x7b73b060 RIP=0x2824737a ntoskrnl.exe!NtClose
CR3=0x7b73b060 RIP=0x28254f0a ntoskrnl.exe!NtHalpSendWaitReceivePort
CR3=0x7b73b1e0 RIP=0x2826b4db ntoskrnl.exe!NtFreeVirtualMemory
CR3=0x7b73b1e0 RIP=0x28268528 ntoskrnl.exe!NtReleaseWorkerFactoryWorker
CR3=0x7b73b1e0 RIP=0x282b6435 ntoskrnl.exe!NtWaitForMultipleObjects
CR3=0x7b73b1e0 RIP=0x28268528 ntoskrnl.exe!NtReleaseWorkerFactoryWorker
CR3=0x7b73b1e0 RIP=0x282685d2 ntoskrnl.exe!NtSetTimer
CR3=0x7b73b1e0 RIP=0x282832c4 ntoskrnl.exe!NtQuerySystemInformation
CR3=0x7b73b1e0 RIP=0x28283964 ntoskrnl.exe!NtQueryInformationProcess
CR3=0x7b73b1e0 RIP=0x282832c4 ntoskrnl.exe!NtQuerySystemInformation
CR3=0x7b73b1e0 RIP=0x28283964 ntoskrnl.exe!NtQueryInformationProcess
CR3=0x7b73b1e0 RIP=0x28283964 ntoskrnl.exe!NtQueryInformationProcess
CR3=0x7b73b1e0 RIP=0x282832c4 ntoskrnl.exe!NtQuerySystemInformation
CR3=0x7b73b1e0 RIP=0x282832c4 ntoskrnl.exe!NtQuerySystemInformation
CR3=0x7b73b1e0 RIP=0x282832c4 ntoskrnl.exe!NtQuerySystemInformation
CR3=0x7b73b1e0 RIP=0x28283964 ntoskrnl.exe!NtQueryInformationProcess
CR3=0x7b73b1e0 RIP=0x28283964 ntoskrnl.exe!NtQueryInformationProcess
CR3=0x185000 RIP=0x2722005 ntoskrnl.exe!ExAllocatePoolWithTag
Heap allocation with known pool tag: 'IoIs' (193497889). nt!Io, I/O SubSystem completion Context Allocation.
CR3=0x7b73b1e0 RIP=0x2826857b1 ntoskrnl.exe!NtWaitForWorkViaWorkerFactory
CR3=0x7b73b1e0 RIP=0x2826857b1 ntoskrnl.exe!NtWaitForWorkViaWorkerFactory
CR3=0x7b73b1e0 RIP=0x2826984b9 ntoskrnl.exe!NtSetTimerEx
CR3=0x7b73b1e0 RIP=0x282832c4 ntoskrnl.exe!NtQuerySystemInformation
CR3=0x7b73b1e0 RIP=0x282832c4 ntoskrnl.exe!NtQuerySystemInformation
CR3=0x7b73b1e0 RIP=0x2826984b9 ntoskrnl.exe!NtSetTimerEx
CR3=0x7b73b1e0 RIP=0x2826984b9 ntoskrnl.exe!NtSetTimerEx
CR3=0x7b73b1e0 RIP=0x2826984b9 ntoskrnl.exe!NtSetTimerEx
CR3=0x7b73b1e0 RIP=0x282b6435 ntoskrnl.exe!NtWaitForMultipleObjects
CR3=0x7b73b1e0 RIP=0x2722005 ntoskrnl.exe!ExAllocatePoolWithTag

```

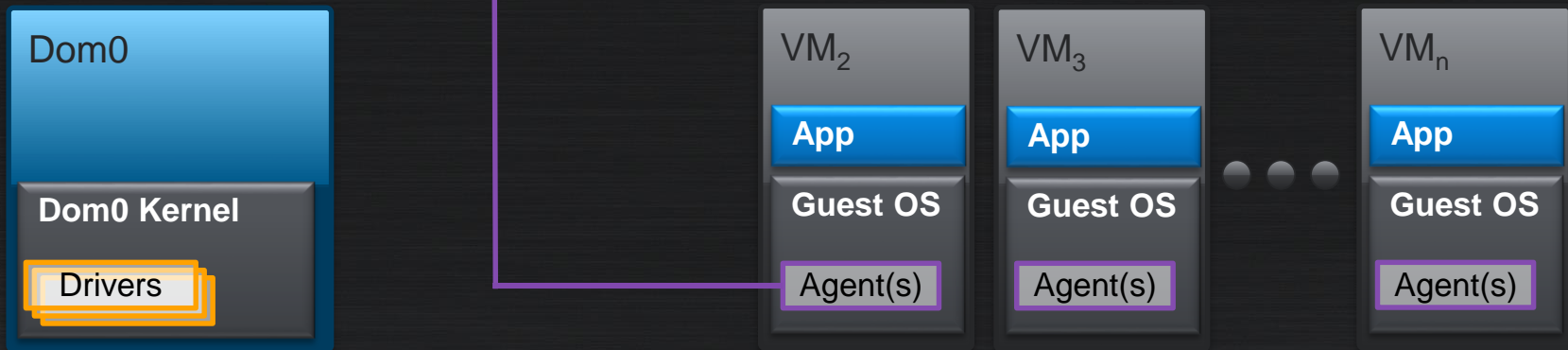


<https://www.youtube.com/watch?v=ZJPHfpDiN4o>

# Cloud Security Today

## Today

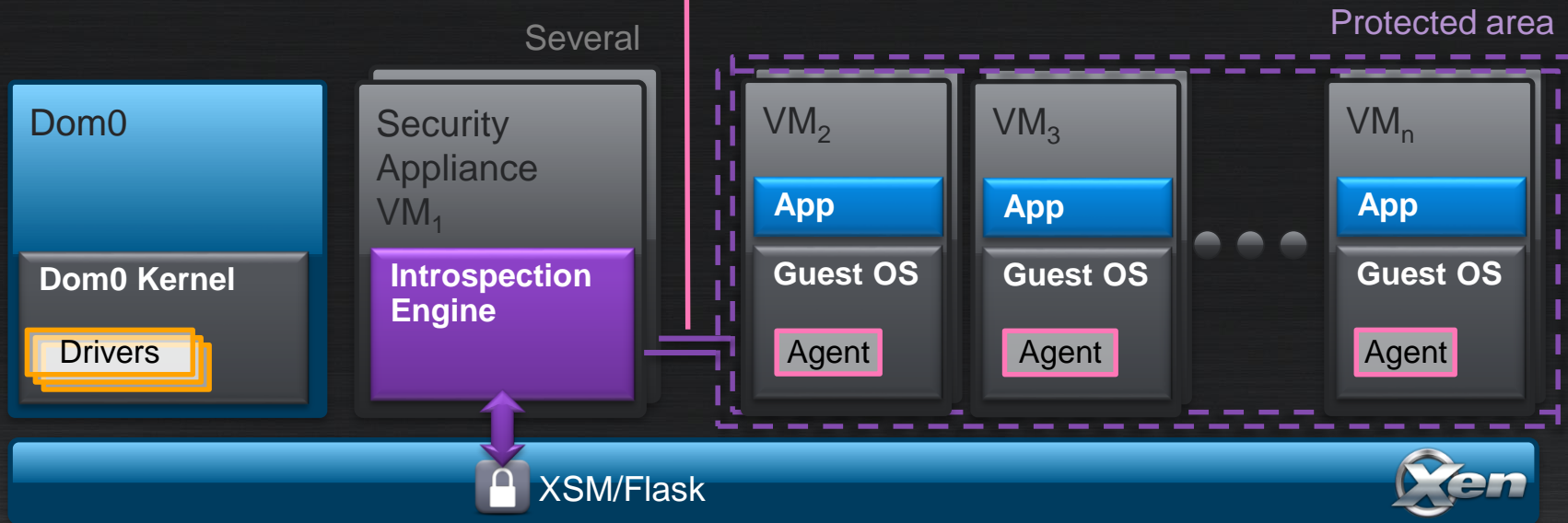
Installed in-guest agents, e.g. anti-virus software, VM disk & memory scanner, network monitor, etc.  
Anti virus storm, deployment/maintenance, ...



# A new model for Cloud Security?

## VMI Approach

Hybrid approach: no need to move everything outside (chose best trade-off)



# Other Security Features being developed

## QEMU and Emulation for Xen secure by default (4.7)

- Response to Venom and other QEMU bugs as an alternative to Stub Domains
- Defense in depth mechanisms to secure the execution of QEMU + Inbuilt Emulation

## Hot Patching or x-splice (4.7+)

- Response to “Cloud Reboots” of 2014 and 2015
- Hypervisor and Workload generation Tooling
- Start with some use-cases and successively add less common ones

## Better Configurability (4.7+)

- Response to criticism from Invisible Things Labs
- Use KCONFIG to disable Hypervisor functionality
- A more wholesome approach to disable and remove undesired functionality

# Conclusion

The project has a history of proactively innovating

The rate of innovation is increasing  
(e.g. more features, more quickly)

The demands on the project are shifting  
(e.g. quality and security, conflicting requirements)

The project has a track record of adapting  
(e.g. to criticism, challenges, ...)

Best Platform for Unikernels in the Cloud  
(e.g. reach, innovation, unique features)

