

Monitoring At Scale: What was Recently Done and What's Next in oVirt

Arik Hadas
Principal Software Engineer
Red Hat
23/10/17

What Do We Mean by “Monitoring”

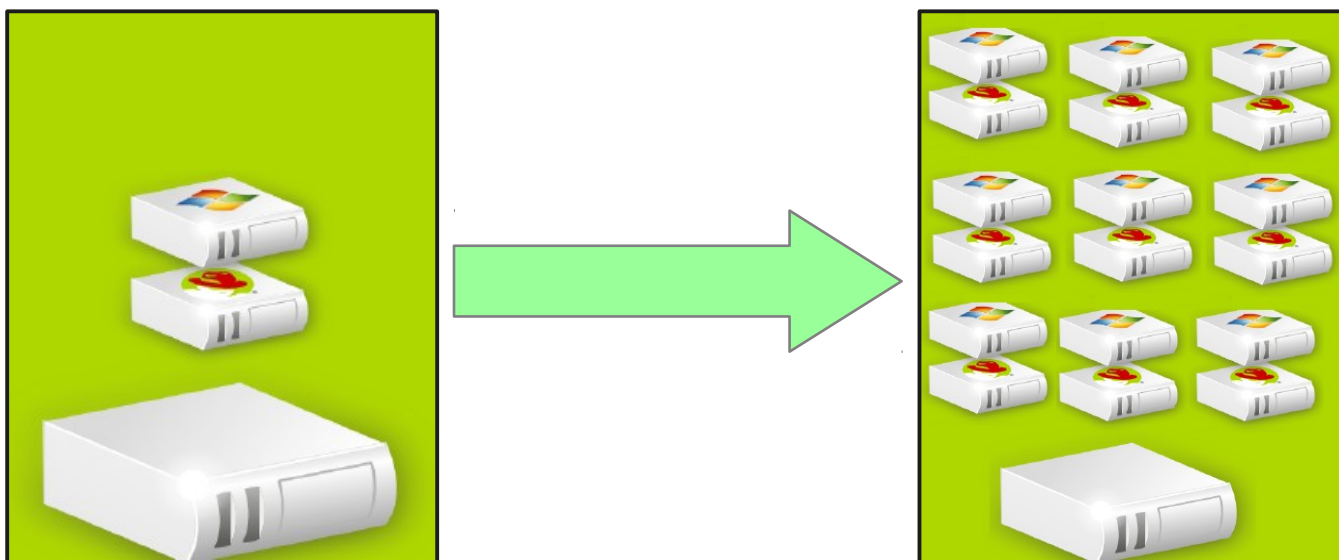
- Identifying the status of active entities
 - VMs, Hosts, Storage domains
- Tracking resource consumption
 - Memory, CPU, Disk space, ...
- Retrieving dynamic properties
 - Client IP, Device addresses, ...

Why is Monitoring Important

- Reflects the up-to-date status of the system
- Affects system responsiveness
- Provides data for automatic processes
 - High availability
 - Load balancing
 - ...

Monitoring At Scale

- The more entities to monitor, the more:
 - Data to collect
 - Data to process
 - Data to store



Problem: Low Performance

- Monitoring
 - Continuous operation
 - Runs in the background
- In large scale deployments monitoring may consume a lot of resources
 - Leads to various anomalies

Our Solution

- Relatively simple changes
 - No architectural change
 - No major change in technology
- We noticed a significant improvement



- **Introduction to oVirt**
- VMs monitoring in large scale deployments
- Improving the monitoring process
- Measurements
- Future work

oVirt What Is oVirt?

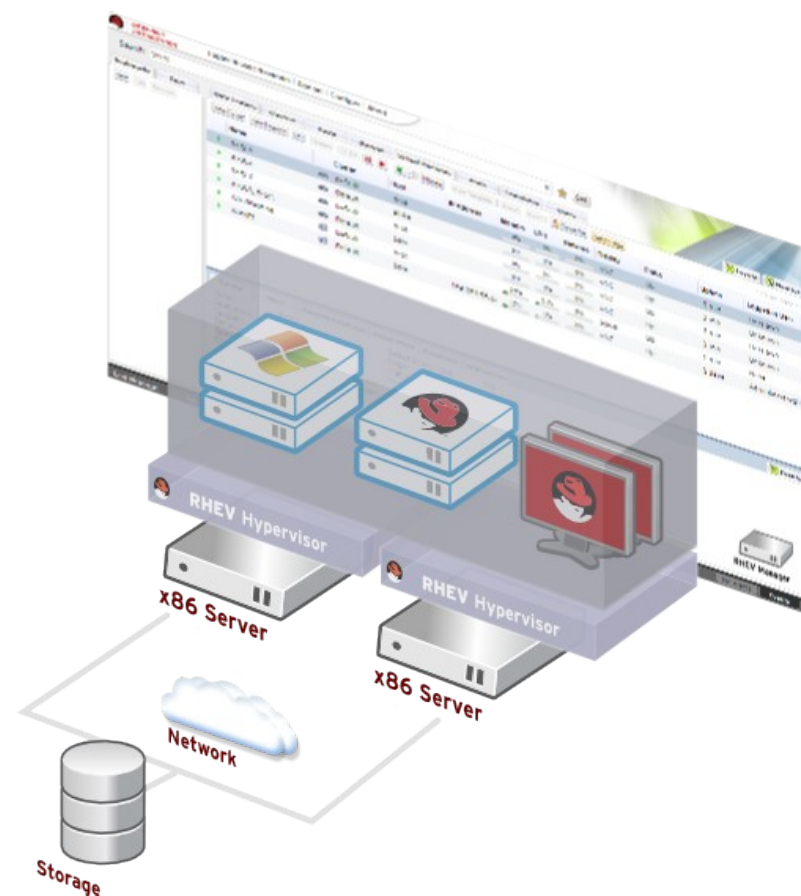
Large scale, centralized management for server and desktop virtualization

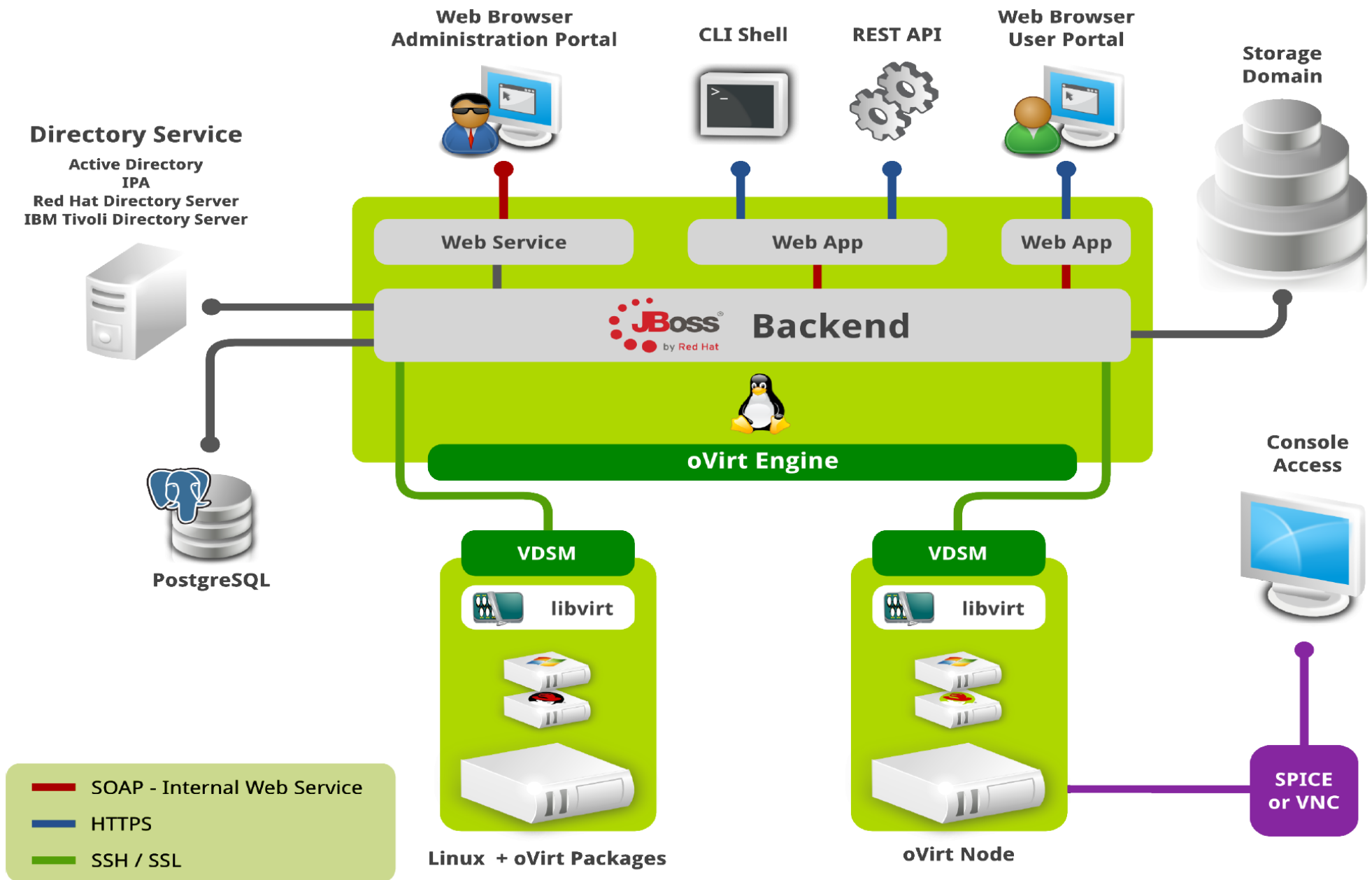
Based on leading performance, scalability and security infrastructure technologies

Provide an open source alternative to vCenter/vSphere

Focus on KVM for best integration/performance

Focus on ease of use/deployment





The screenshot displays the oVirt Webadmin interface. The top navigation bar includes the oVirt logo, the text "OPEN VIRTUALIZATION MANAGER", and several utility icons. A left-hand sidebar contains a menu with categories like Dashboard, Compute, Network, Storage, Administration, and Events. The "Compute" section is expanded, showing sub-items: Virtual Machines (selected), Templates, Pools, Hosts, Data Centers, and Clusters. The main content area features a toolbar with actions such as "New VM", "Edit", "Remove", "Run", "Suspend", "Shutdown", "Reboot", "Console", "Migrate", and "Create Snapshot". Below the toolbar is a table listing virtual machines. The table has columns for Comment, Host, IP Address, FQDN, Cluster, Data Center, Memory, CPU, Network, Graphics, Status, Uptime, and Description. Three VMs are listed, all associated with "test-cluster" and "test-dc". The first VM is "lago-basic-suite-...", which is "Up" with 1 hour of uptime. The other two VMs are "Down".

Comment	Host	IP Address	FQDN	Cluster	Data Center	Memory	CPU	Network	Graphics	Status	Uptime	Description
	lago-basic-suite-...			test-cluster	test-dc	0%	1%	0%	SPICE + V...	Up	1 h	
				test-cluster	test-dc	--	--	--	None	Down		CirROS imported from ...
				test-cluster	test-dc	--	--	--	None	Down		

https://192.168.201.4/ovirt-engine/webadmin/?locale=en_US#vms

Feature-Rich Platform

VM
Affinity

Hosted
Engine

Hostdev
Passthrough

Live
Snapshots

Neutron
Integration

Port
Mirroring

Foreman
Integration

VM High
Availability

Glance
Integration

Live
Storage
Migration

Multi-Level
Administration

V2V

Backup
API

Gluster
Support

Live
Merge

VM
Migration

Cloud-init
Integration

Hot Plug
Memory

Nvidia
GPU

VM
Pools

Direct
LUNs

VM
Leases

Load
Balancing

Hot Plug
CPU

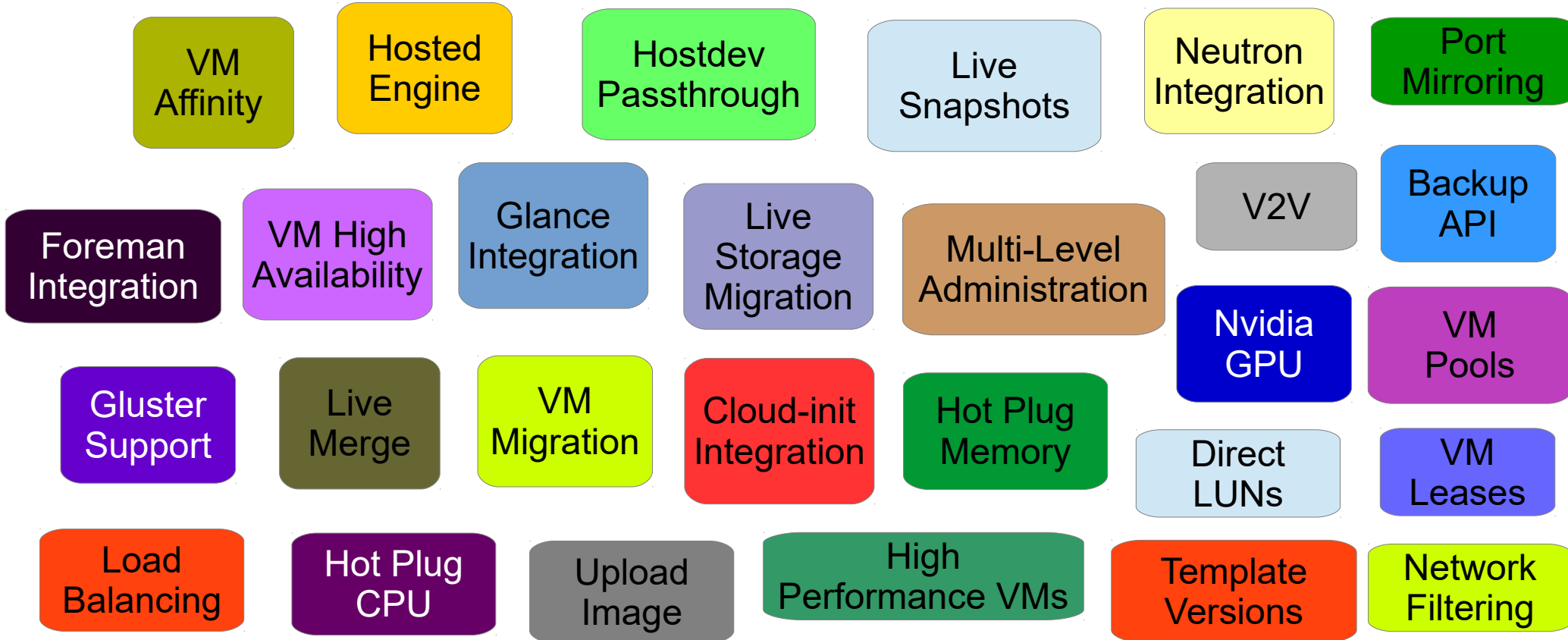
Upload
Image

High
Performance VMs

Template
Versions

Network
Filtering

Feature-Rich Platform



- Less attention to scale

- Introduction to oVirt
- **VMs monitoring in large scale deployments**
- Improving the monitoring process
- Measurements
- Future work

- Focus on monitoring of virtual machines
 - Far more instances than any other entity
- This includes:
 - Status
 - Dynamic properties (i.e., client IP)
 - Devices information
 - Statistics

-
- The diagram illustrates the integration of JBoss Backend with VDSM and libvirt. On the left, a green box represents the **oVirt Engine**, which contains a **JBoss Backend** (JBoss by Red Hat) and a Tux penguin icon. Above the JBoss Backend are three components: **Web Service**, **Web App**, and **Web App**. On the right, a green box represents the **VDSM** (Virtual Desktop Service Manager) layer, which contains **libvirt** and three server icons. Two large green curved arrows indicate the flow of data and control between the JBoss Backend and the VDSM/libvirt layer.

- Hosts are locked during monitoring cycles
 - To prevent operations on VMs in parallel
- Dynamic properties are compared via reflection
- VM statistics are not being compared
 - They almost always change
- Devices are polled separately when their hash changes

- Problems in very large scale deployments
 - Monitoring cycles were skipped
 - High CPU consumption
 - High load on the database

Proposed Solutions

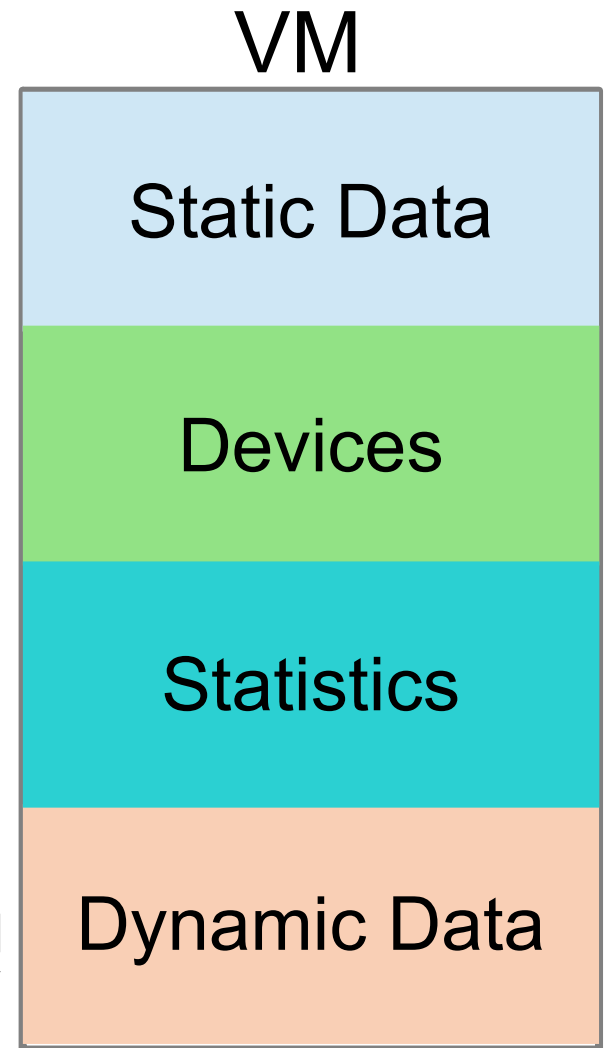
- Add a global caching layer
 - To reduce interactions with the database
 - Does not solve the high CPU consumption
- Distribute the monitoring process
 - Addresses the high CPU consumption
 - Does not reduce the load on the database
- Both solutions were too complex

- Introduction to oVirt
- VMs monitoring in large scale deployments
- **Improving the monitoring process**
- Measurements
- Future work

oVirt Issue #1: Too Many Writes to DB

- Static data is not monitored
- Devices rarely change
- Statistics change in each cycle
- Some of the dynamic data (reported data) might change
 - Not often though

Reported data (i.e., client IP) +
Not reported data (i.e., stop reason)



Reduce Number of Writes

- Introduce @UnchangeableByVdsm
 - Marks properties that are not reported

```
private String currentCd;  
@UnchangeableByVdsm  
private String stopReason;  
private VmExitReason exitReason;
```

- Move frequently changed fields to the stats
 - E.g., guest memory cached/buffered/free

- Devices hash was stored with the dynamic data
 - Consequently, change of one device triggered persistency of all dynamic data
- Solution: store the devices hash separately

Issue #2: Too Many Reads from DB

- Many connections with DB are used
- Long time is spent on querying the DB
- Even when no data (except stats) is changed!

Eliminate Redundant Queries

- Optimize the code to skip unneeded data processing (including queries from DB)
- For example, skipping redundant VM numa nodes processing eliminated the following DB interactions:

Average time (micro-sec)

- 261 to get numa nodes by host
- 259 to get assigned numa nodes
- 255 to get numa node CPU by host
- 246 to get numa node CPU by VM
- 242 to get numa nodes by VM

Overall time (micro-sec)

- Getting numa nodes by host—3% (48,546 msec)
- Getting assigned numa nodes—3% (48,201 msec)
- Getting numa node CPU by host—3% (47,569 msec)
- Getting numa node CPU by VM—2% (45,918 msec)
- Getting numa nodes by VM—2% (45,041 msec)

- Apply memoization to repeated queries

```
public class MemoizingSupplier<T> implements Supplier<T> {  
    private final Supplier<T> delegate;  
    private boolean initialized;  
    private T value;  
  
    public MemoizingSupplier(Supplier<T> delegate) {  
        this.delegate = delegate;  
    }  
  
    public T get() {  
        if (!initialized) {  
            value = delegate.get();  
            initialized = true;  
        }  
        return value;  
    }  
}
```

- Cache only relevant entity's properties
 - E.g., static properties used by the monitoring
- Cache only relevant entities
 - E.g., VM jobs (limited number of instances)
- Use DB for persistency, not as a bus of data
 - E.g., VM statistics

Lighter, Dedicated Queries

- Complicated queries take time
- Attempt #1: narrow down 'vms' view

> explain analyze select * from vms where ...

Planning time: 2.947 ms

Execution time: 765.774 ms

> explain analyze select * from vms_monitoring_view where ...

Planning time: 0.387 ms

Execution time: 275.600 ms

- Attempt #2: query only dynamic data

> explain analyze select * from vms_monitoring_view where ...

Planning time: 0.405 ms

Execution time: 275.850 ms

> explain analyze select * from vm_dynamic where ...

Planning time: 0.109 ms

Execution time: 2.703 ms

Issue #3: Locks Contention

- High contention between monitoring threads and those executing operations on VMs
- During the execution of VM operations, the host was locked to avoid monitoring the VM
 - To prevent conflicts

- Replaced host-level locks with VM-level locks
 - VM operations lock VMs rather than hosts
 - Monitoring locks each VM running on the host
 - And skips those that cannot be locked
- That reduces contention rate on operations-intensive deployments

Issue #4: High UNIX Load

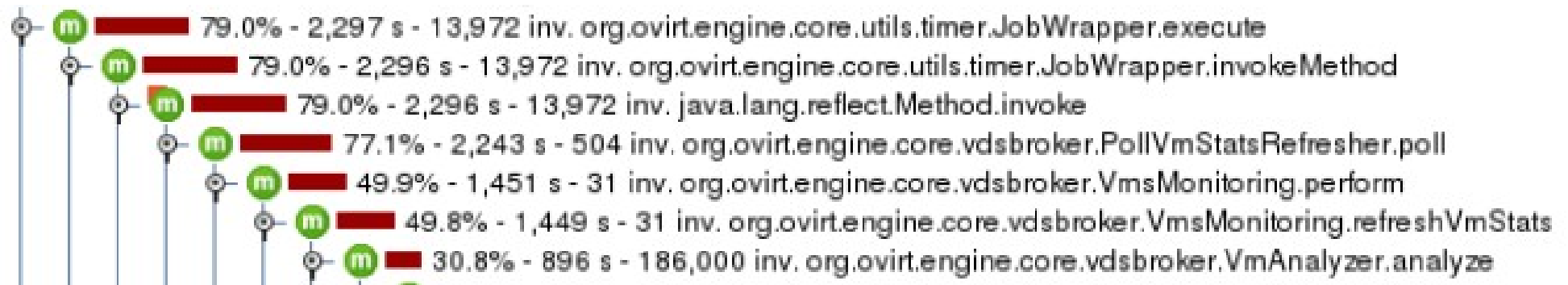
- The overall backend load was relatively high
 - Even in stable deployment
- The monitoring was an immediate suspect

- Replaced the polling-based backend \leftrightarrow host protocol with events-based protocol
 - Based on JSON-RPC instead of XML-RPC
- Hosts send events upon VM changes
 - Less monitoring cycles and data to process
- Keep polling statistics cycles
 - Statistics always change
 - Compensate missing events

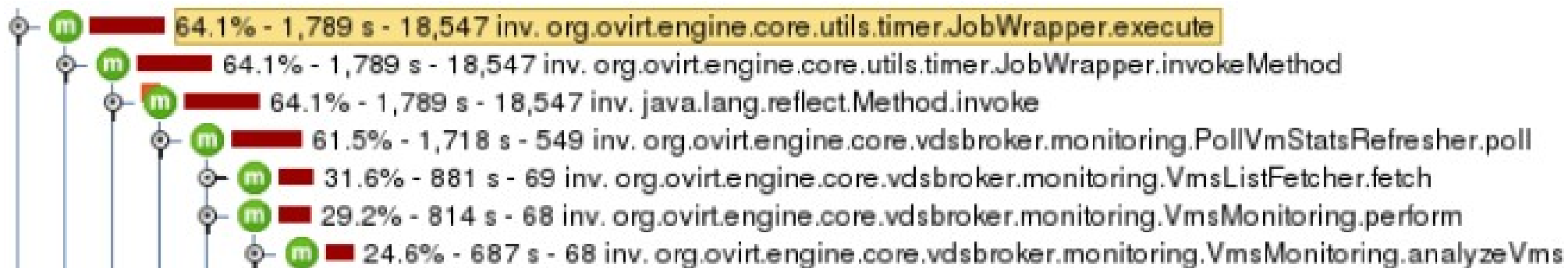
- Introduction to oVirt
- VMs monitoring in large scale deployments
- Improving the monitoring process
- **Measurements**
- Future work

- Deployment with 1 host running 6000 VMs
 - ‘Fake VMs’
- Stable deployment
 - No operation is done
- Measured 1 hour of uptime
- Compared versions 3.6 and 4.1
 - Both used events

3.6

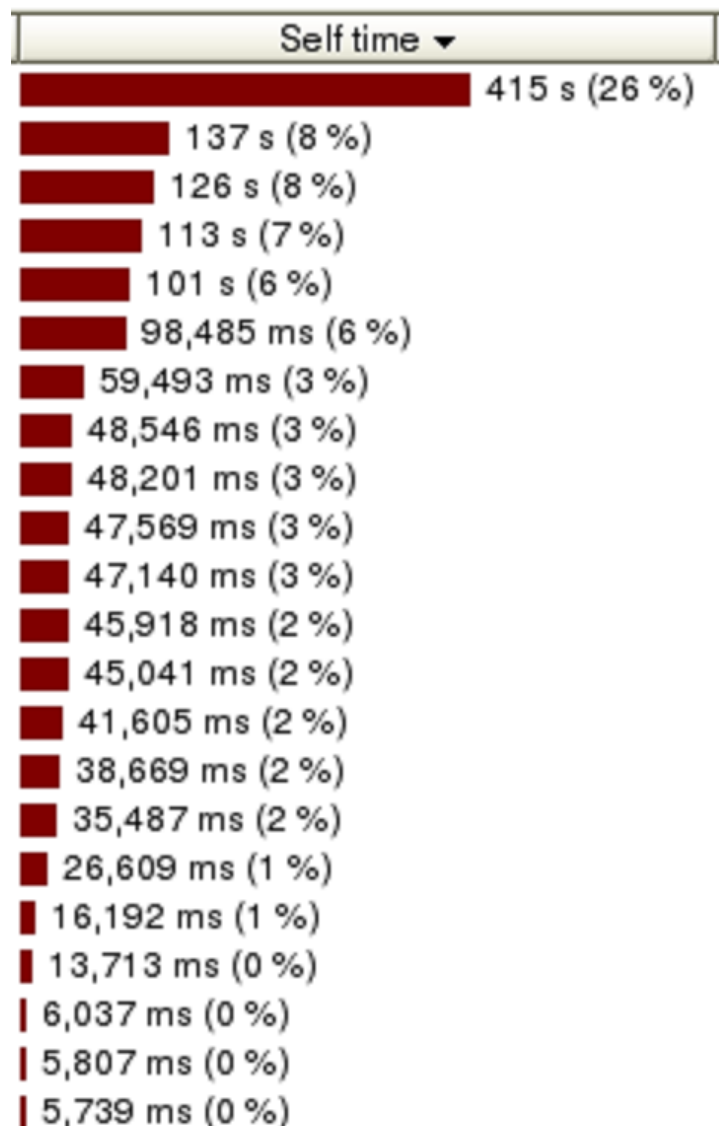


4.1

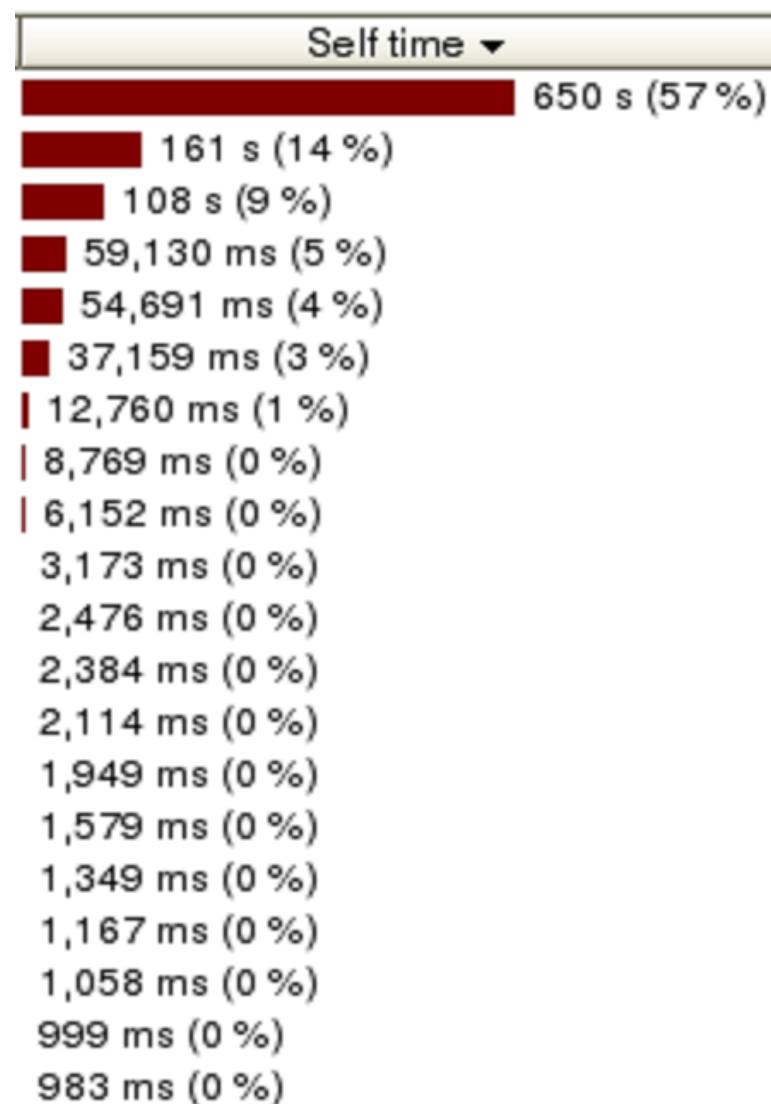


- Total CPU time reduced from 2297s to 1789s (78%)
- Significantly less time in monitoring code
 - Processing time reduced from 896s to 687s
 - Persistence time reduced from 546s to 114s
 - Overall, 814s instead of 1451s (56%)

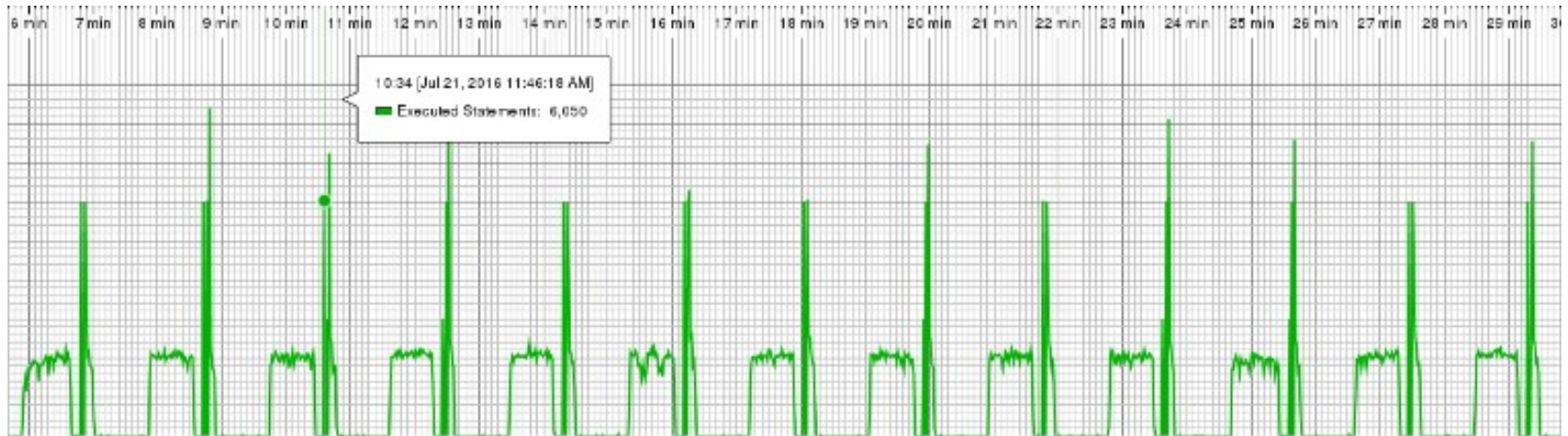
3.6



4.1



3.6

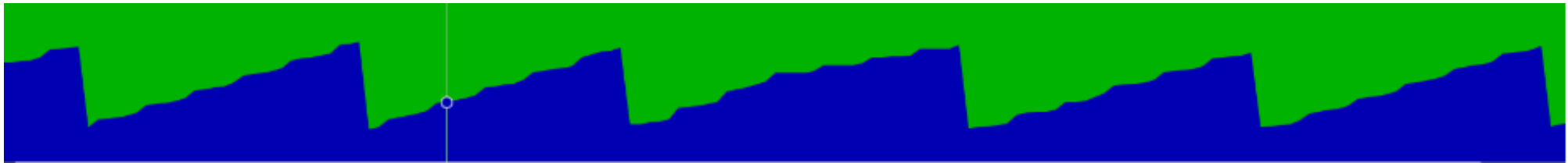


4.1

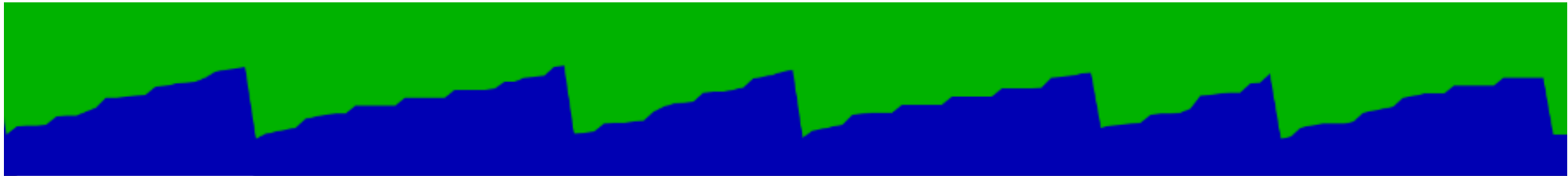


- The time to query all VMs reduced from 3539ms to 909msec (26%)
- The time to save dynamic data in 3.6 was 101 sec (6%, 544 micro-sec on average), 0 in 4.1
 - Similar results for other properties
- In overall, less use of the database

3.6



4.1



- Surprisingly, less memory was consumed in 4.1
 - In 3.6 it gets to ~1.45GB
 - In 4.1 it gets to ~1.2GB
- Probably because of caching done by postgres

- Introduction to oVirt
- VMs monitoring in large scale deployments
- Improving the monitoring process
- Measurements
- **Future work**

- Separate out statistics monitoring
- Apply similar principles to host monitoring
- Add caching of more entities
 - Specifically, VM dynamic data (e.g., status)

- Significant improvement shown in a case study
 - All changes are available in version 4.1
- This required deep knowledge of the platform
 - No shortcuts in the form of generic solutions
 - No major technological change
 - No architectural change

THANK YOU!

<http://www.ovirt.org>
ahadas@redhat.com
[ahadas@irc.oftc.net#ovirt](irc://irc.oftc.net/#ovirt)