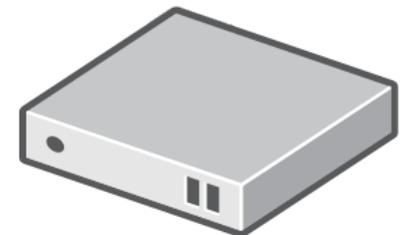
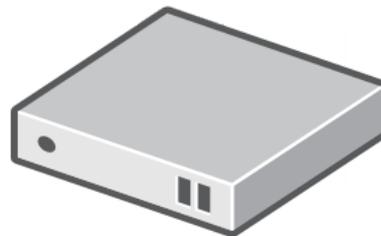
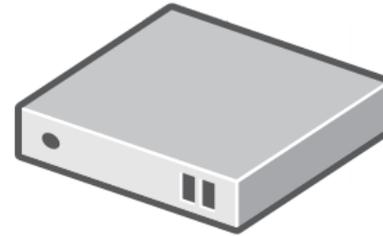
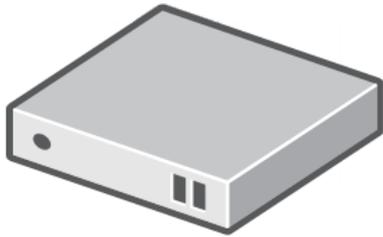


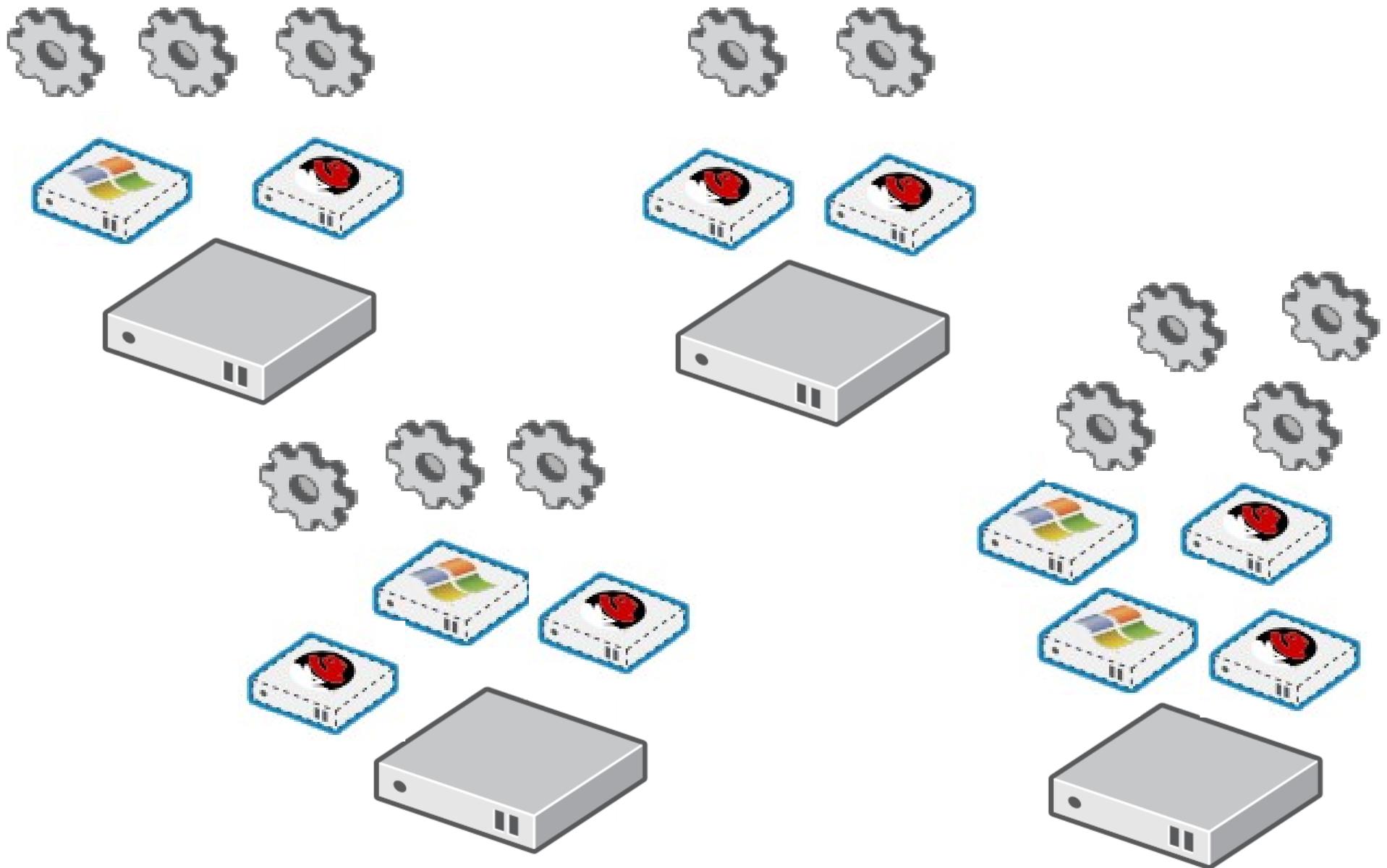
# High Availability with No Split Brains!

Arik Hadas  
Principal Software Engineer  
Red Hat  
27/01/2018

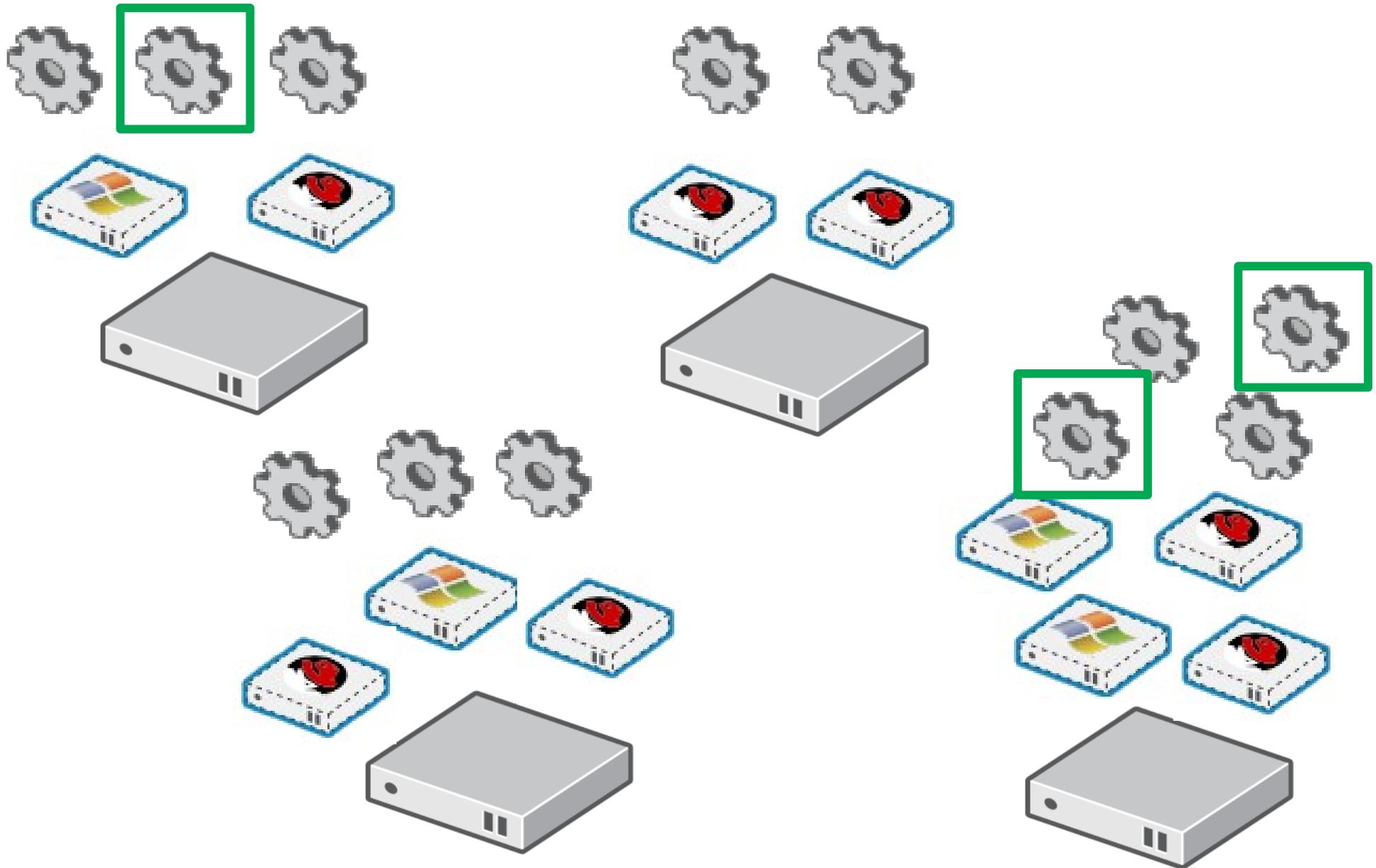


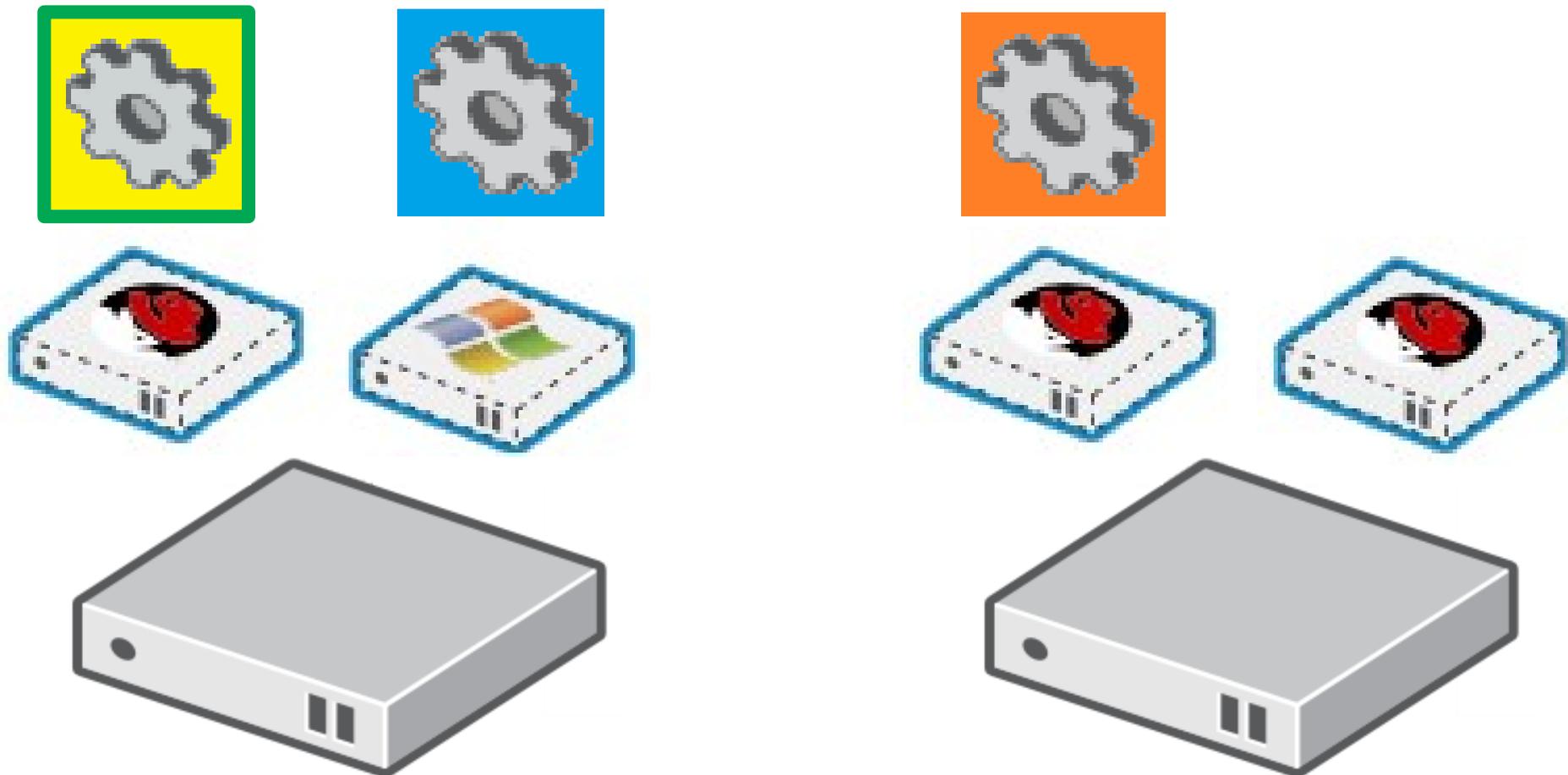


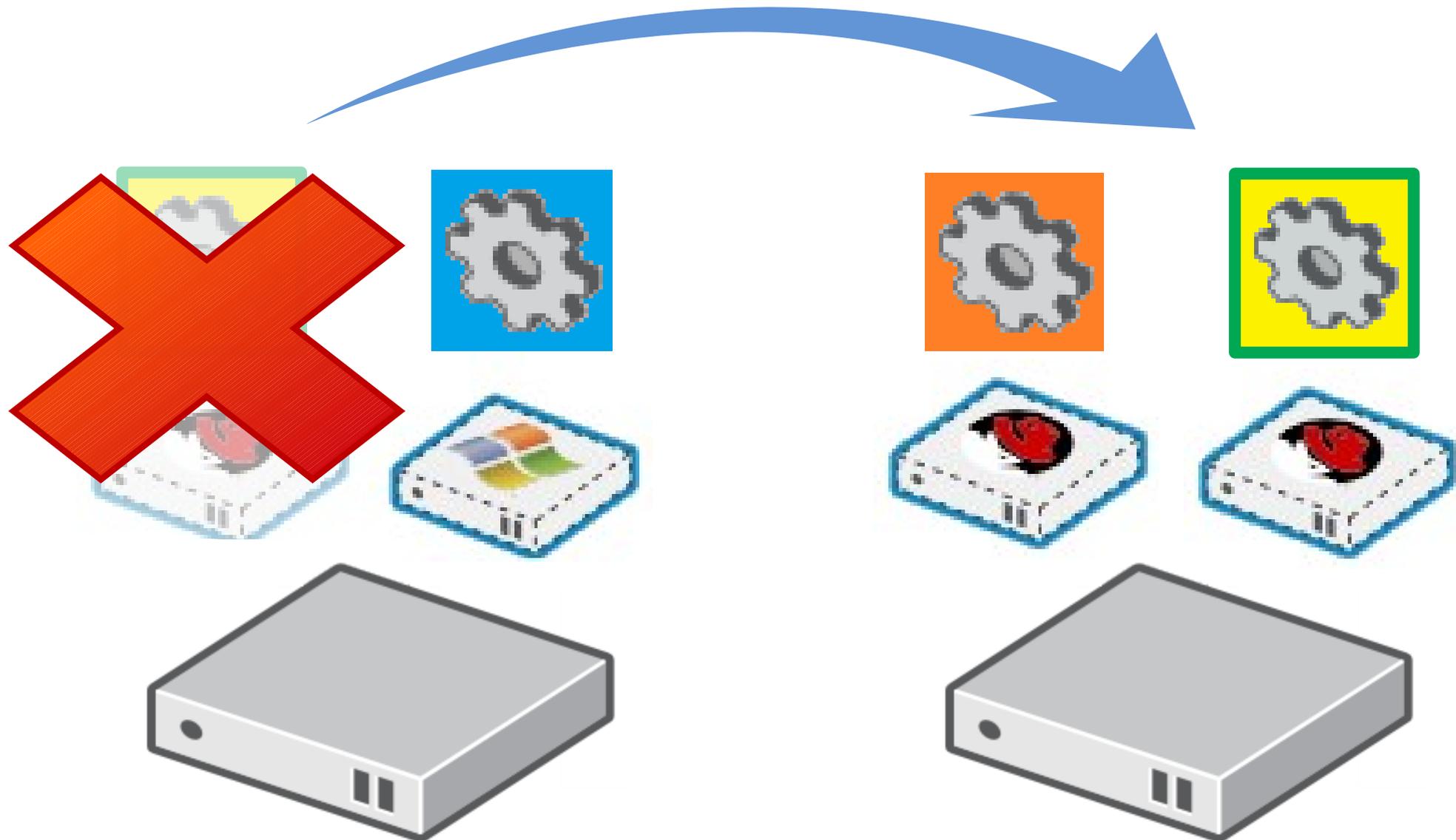
# oVirt Virtual Data Center - Applications



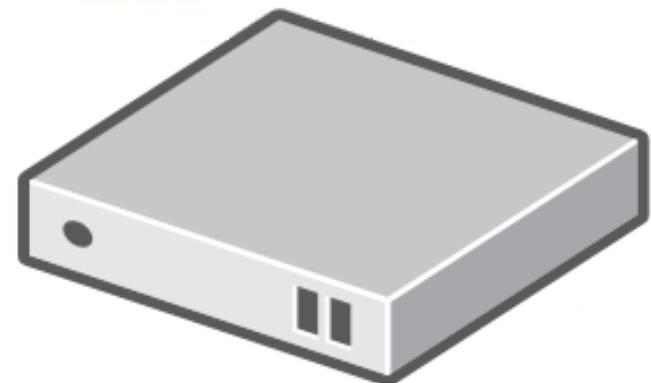
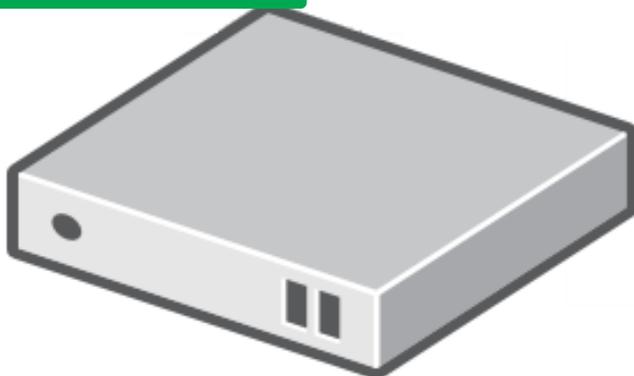
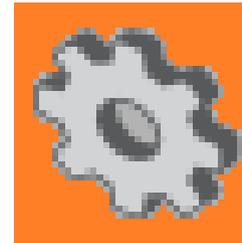
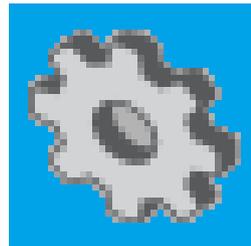
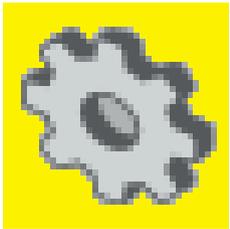
# oVirt Some Applications are More Critical



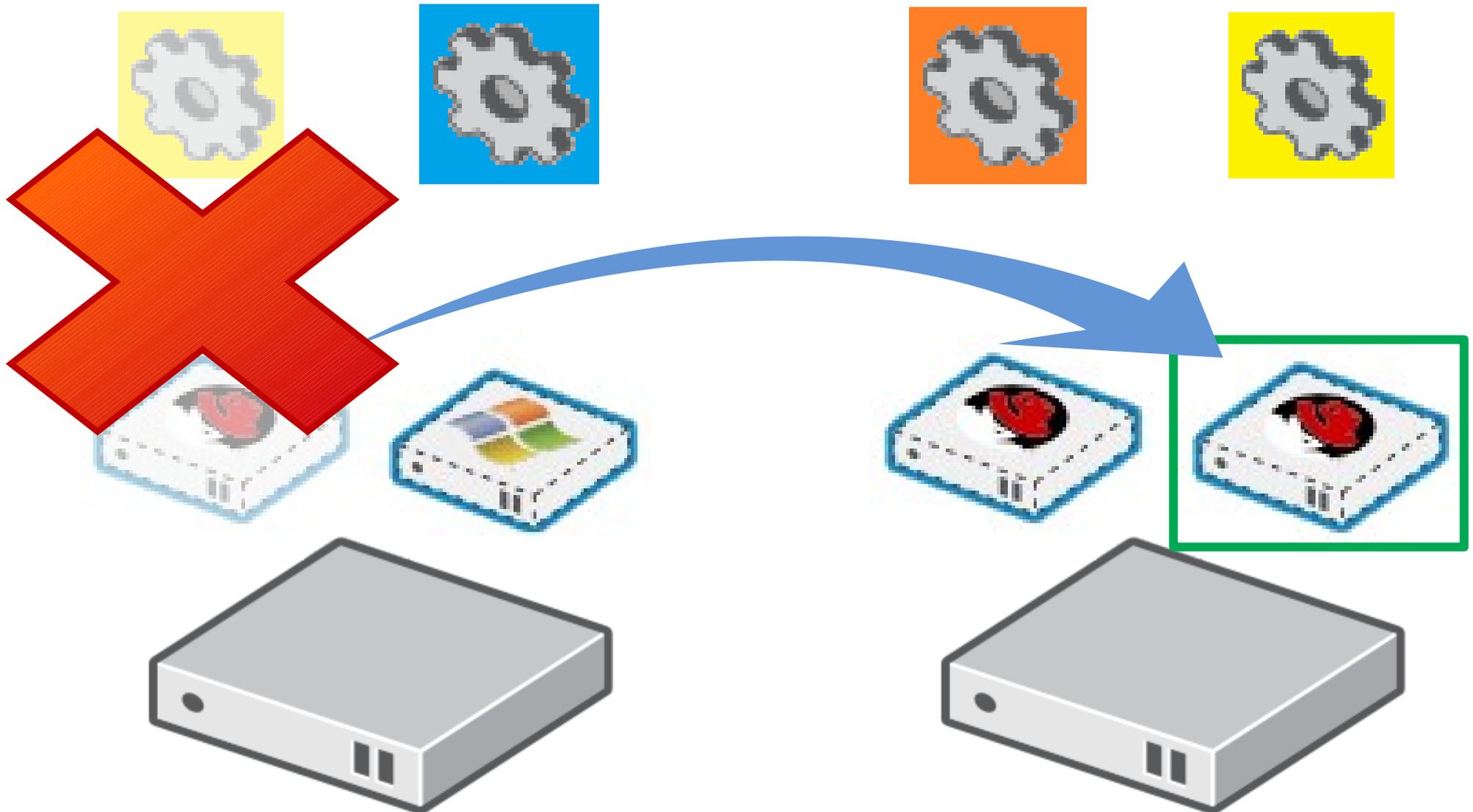




- Higher resource consumption
- More responsibility on the application
- Backup starts in a different environment
  - Different IP address(es)
  - Different disk(s)

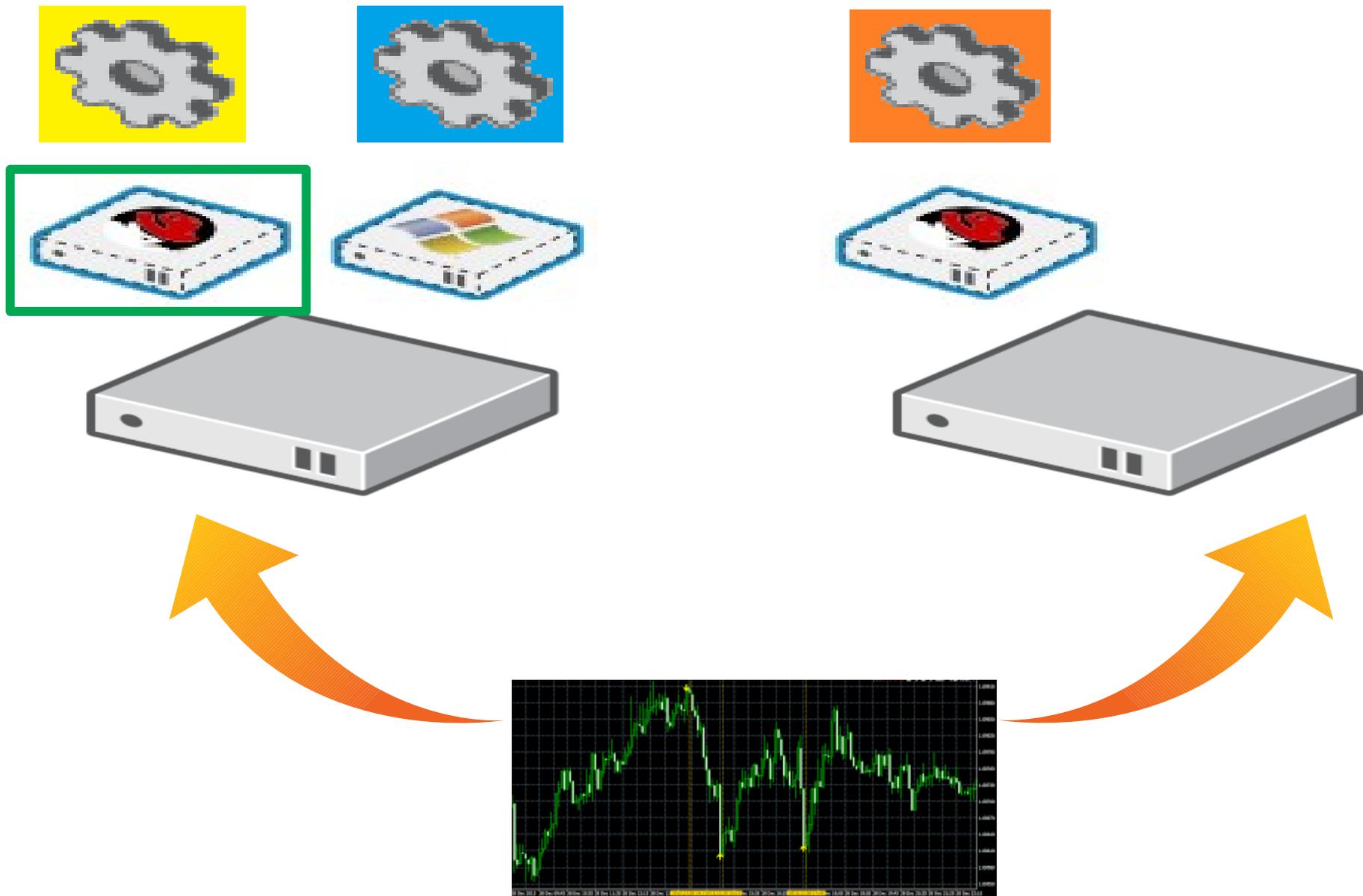


# High Availability - VM-Level

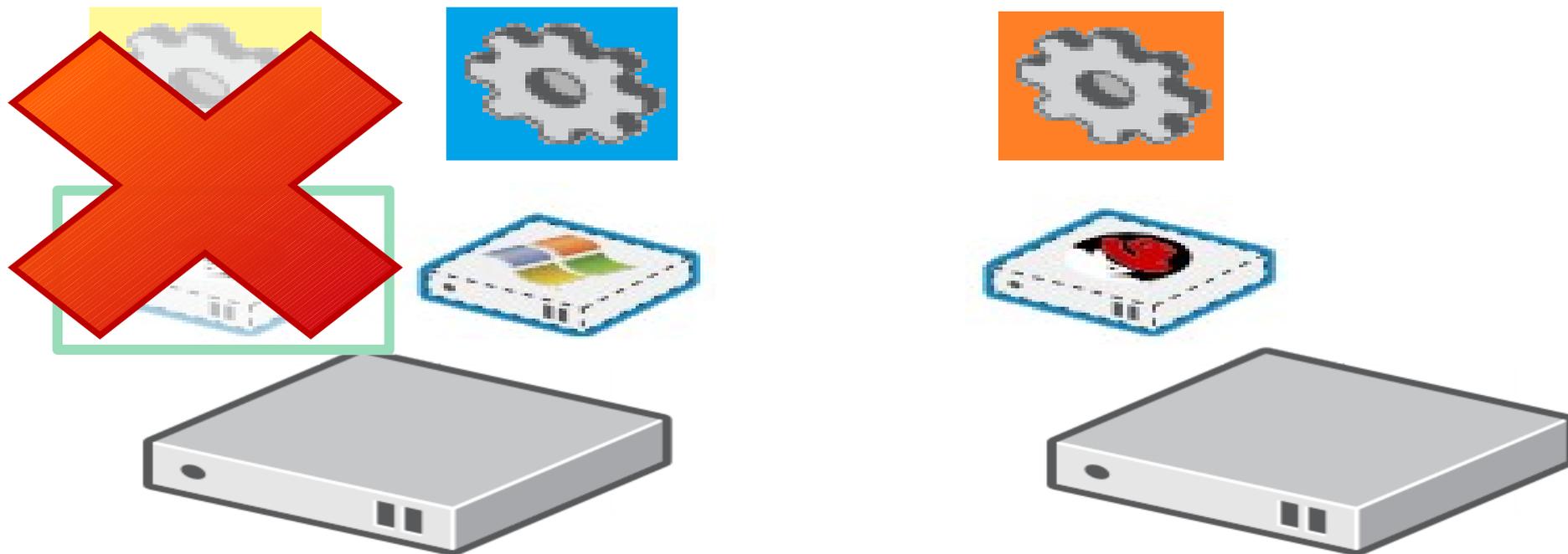


- More efficient resource consumption
- Implemented at the infrastructure level
- VM always start in the same environment
  - Same IP address(es)
  - Same disk(s)

# Central Monitoring Unit



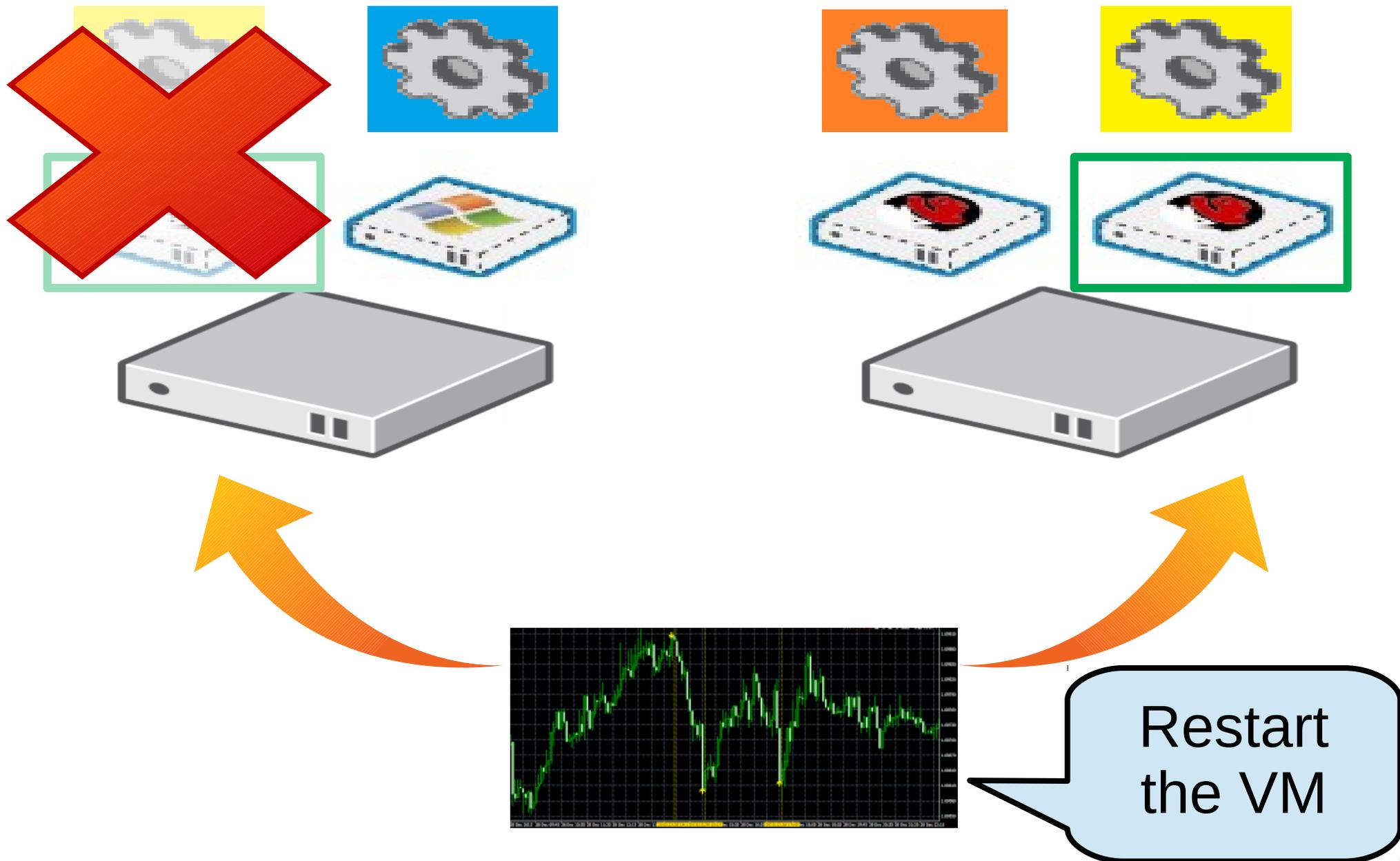
# oVirt Fault Detection



HA VM  
went down!

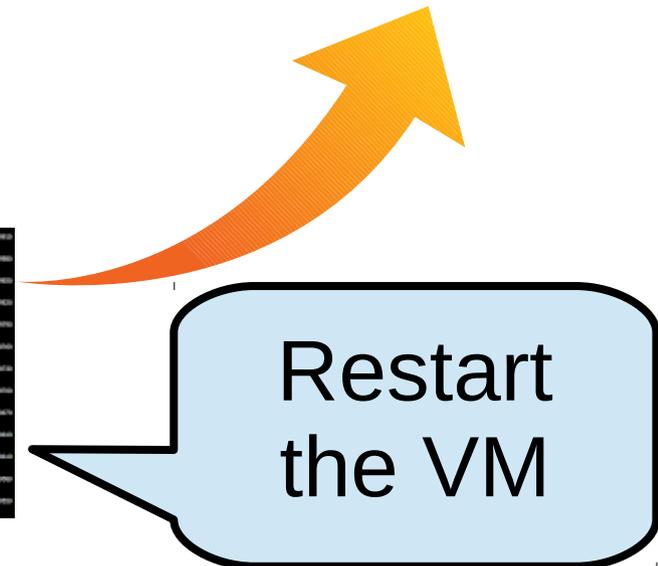
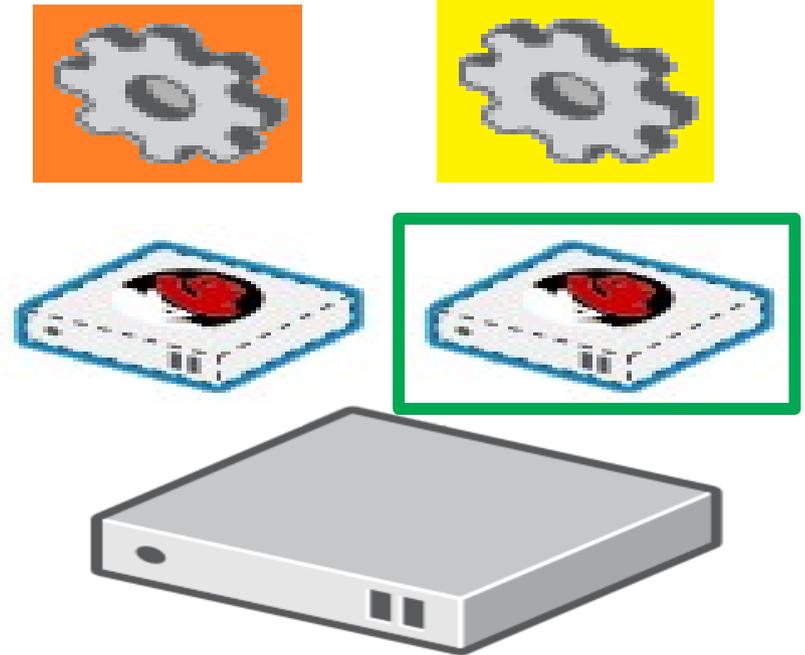


# oVirt Automatic Restart



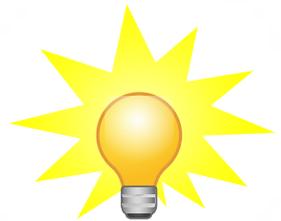
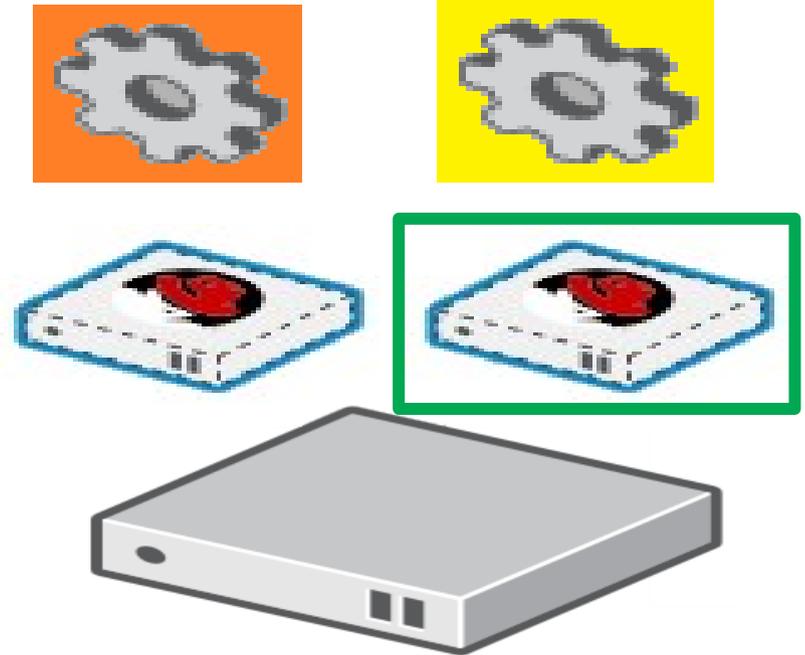
What if:

- Inaccessible resources
- VM is locked
- VM is being intentionally shut down



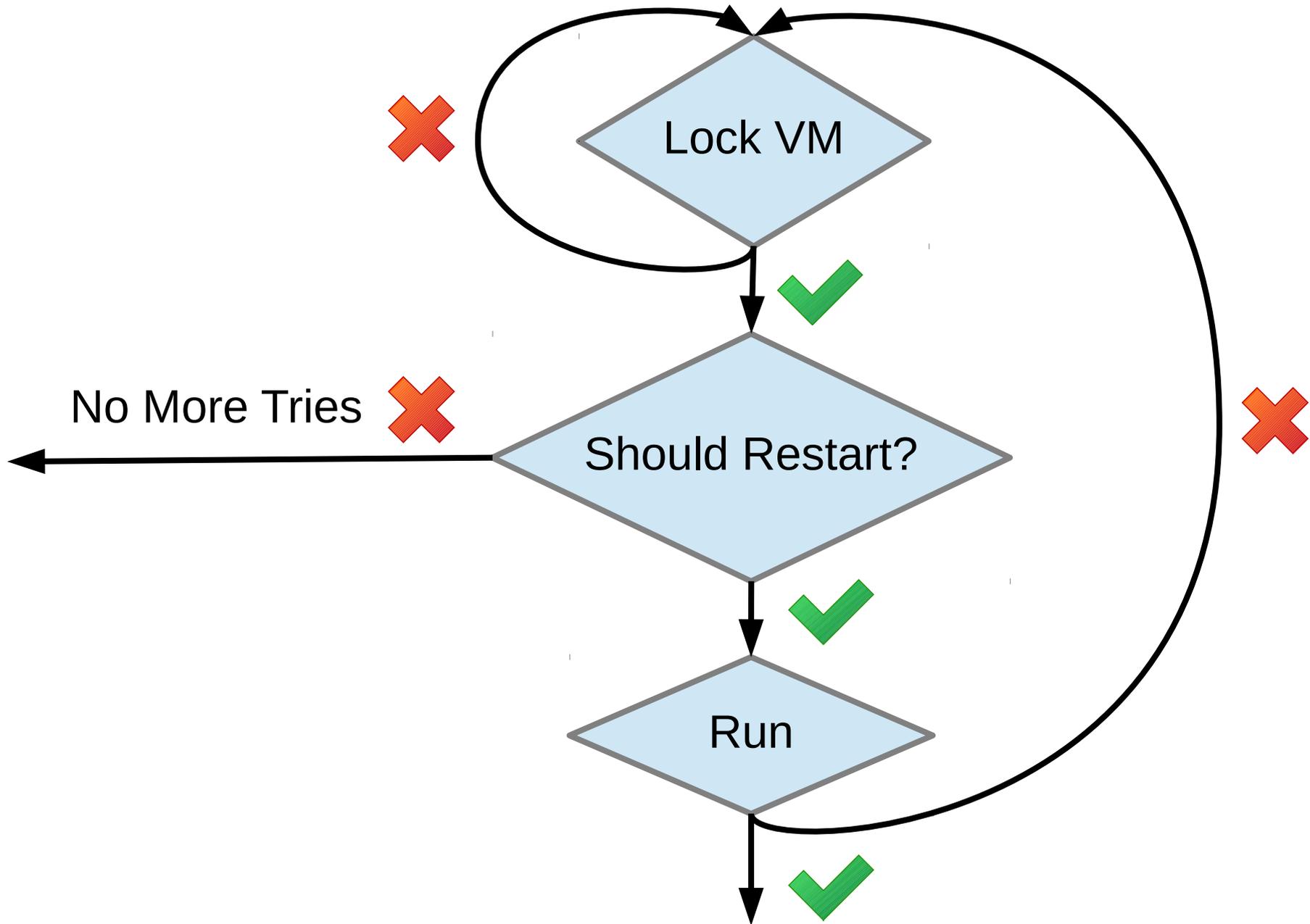
What if:

- Inaccessible resources
- VM is locked
- VM is being intentionally shut down

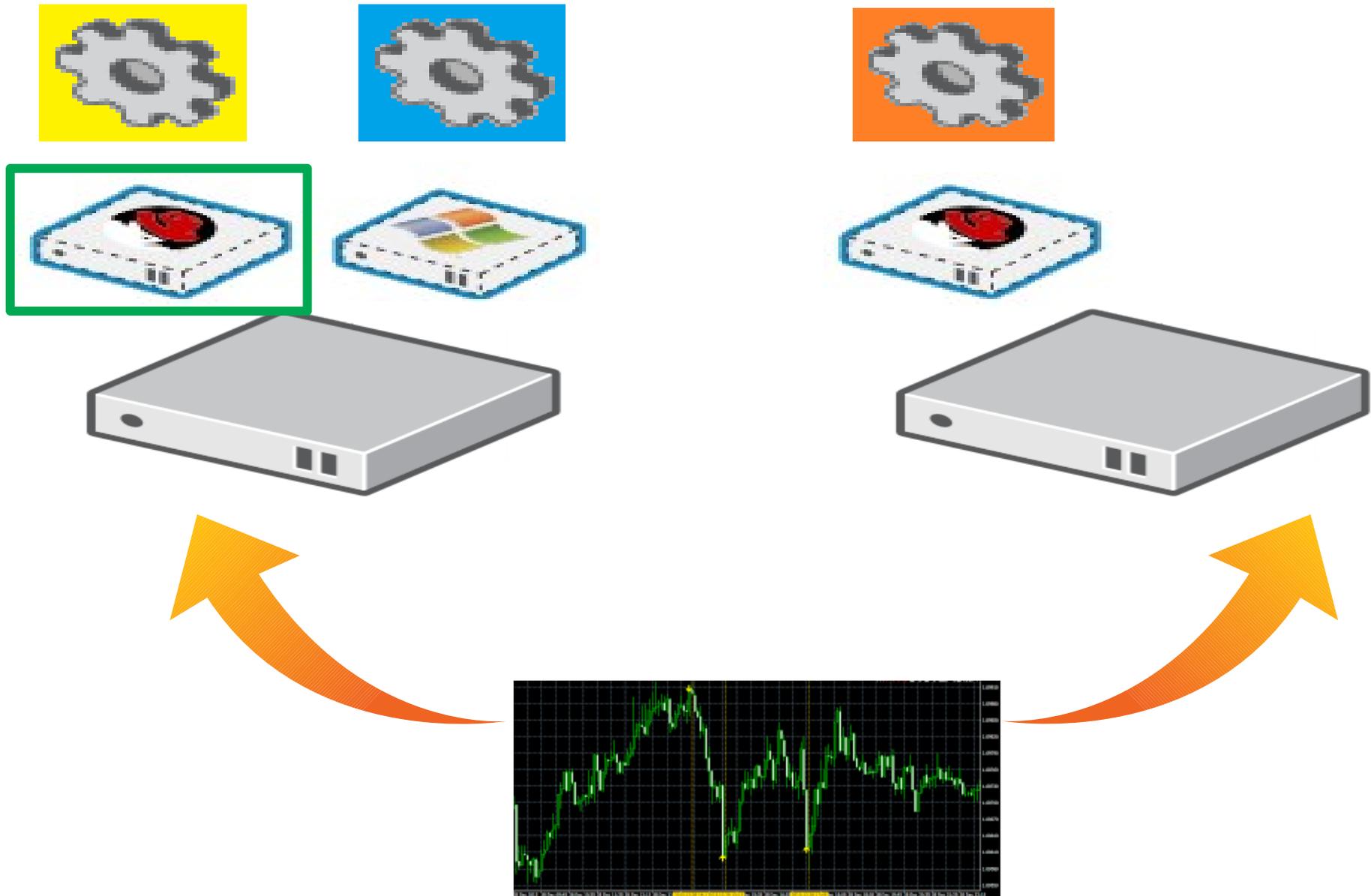


## AutoStartVmsRunner

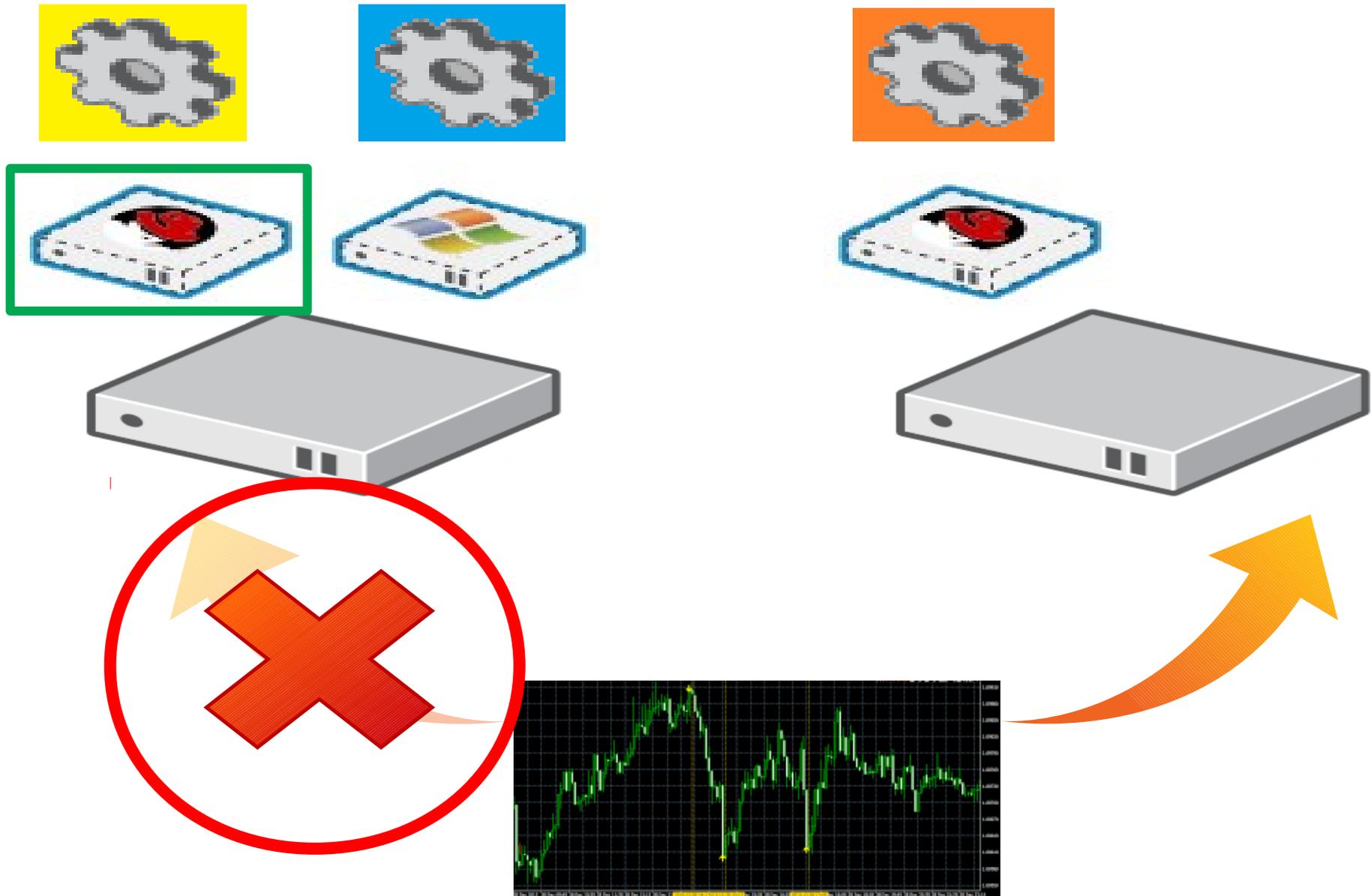
<https://github.com/oVirt/ovirt-engine/blob/master/backend/manager/modules/bll/src/main/java/org/ovirt/engine/core/bll/AutoStartVmsRunner.java>



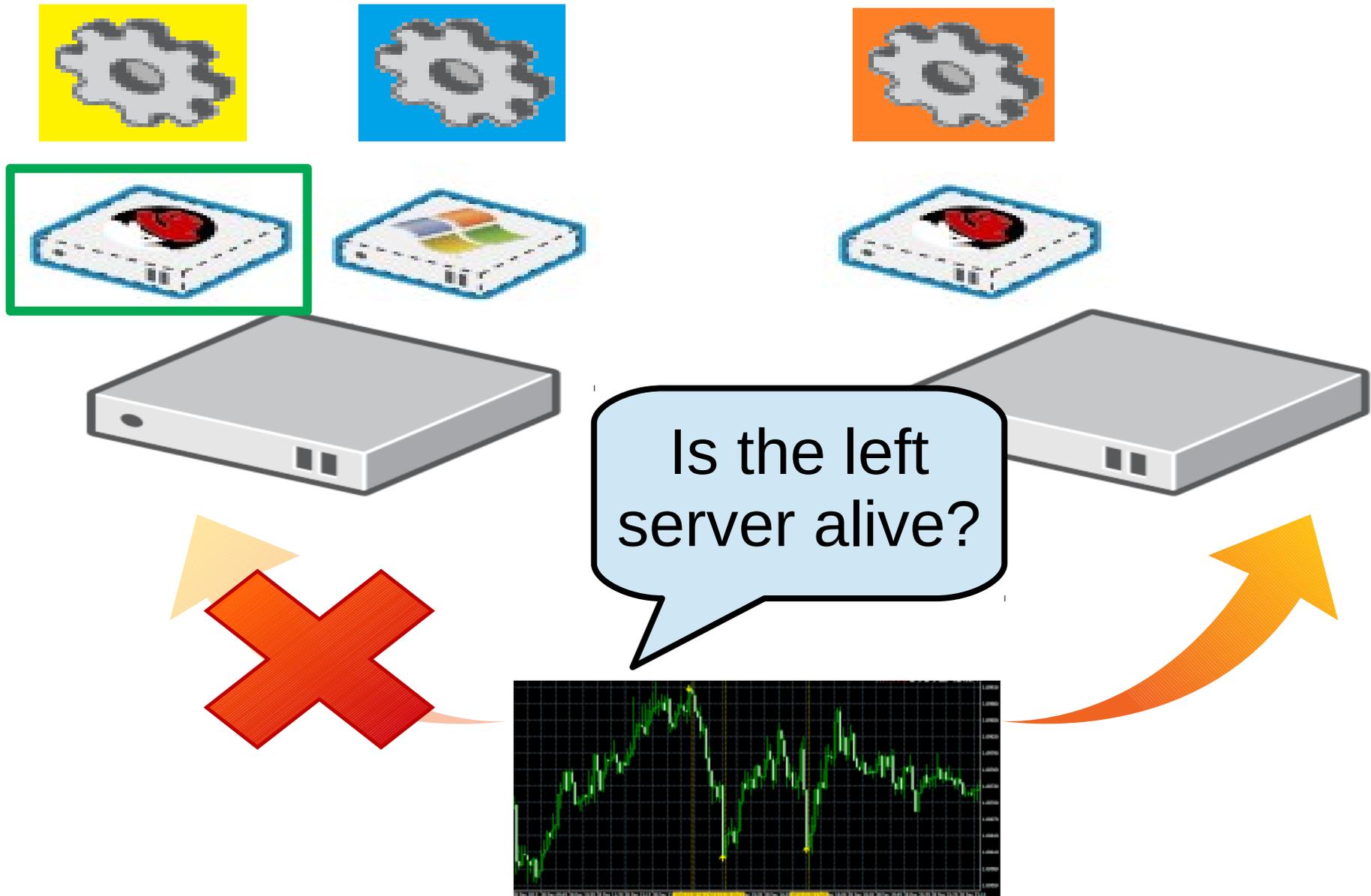
# oVirt Fault Detection - Even More Complex



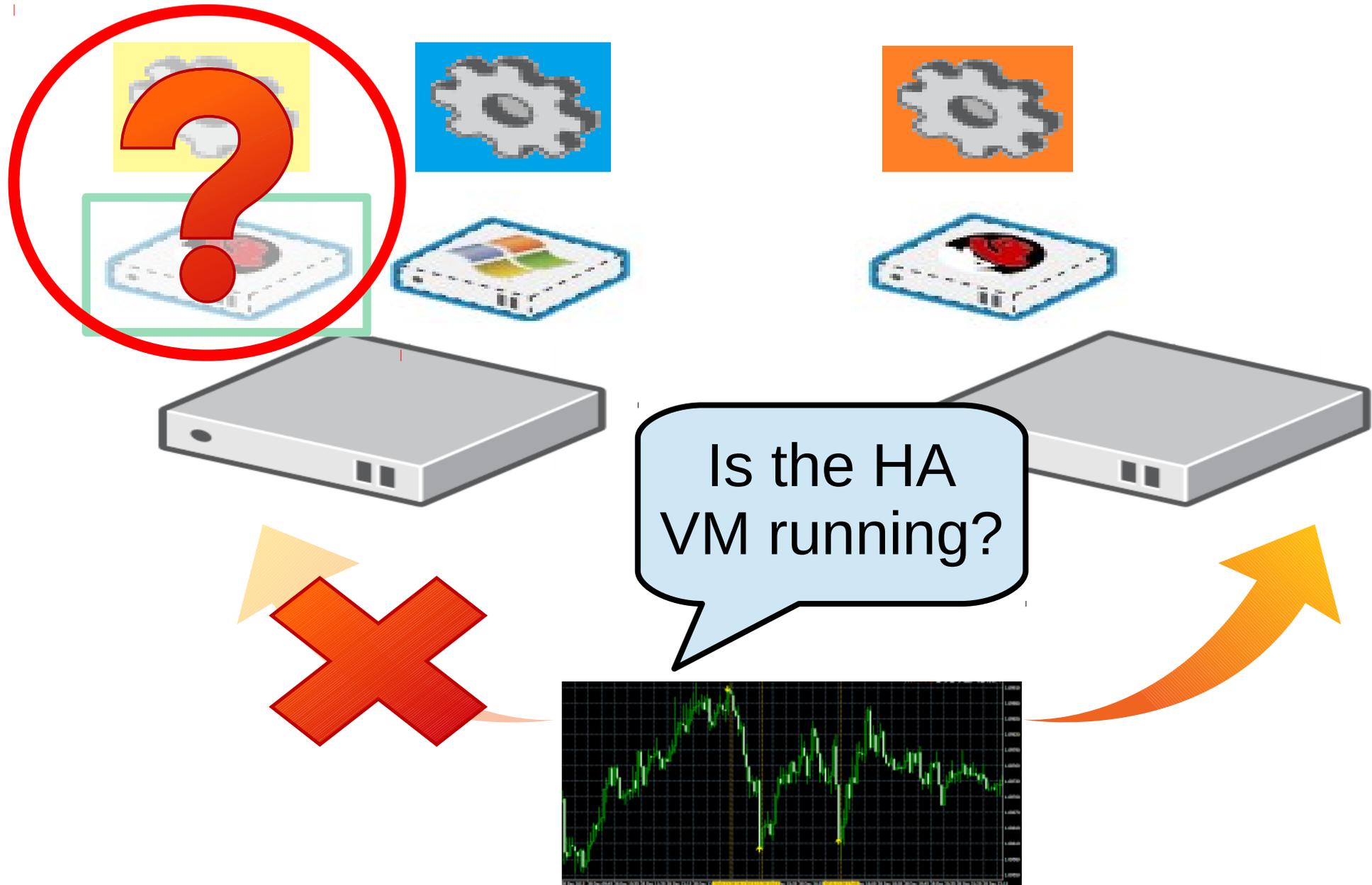
# oVirt Fault Detection – Even More Complex



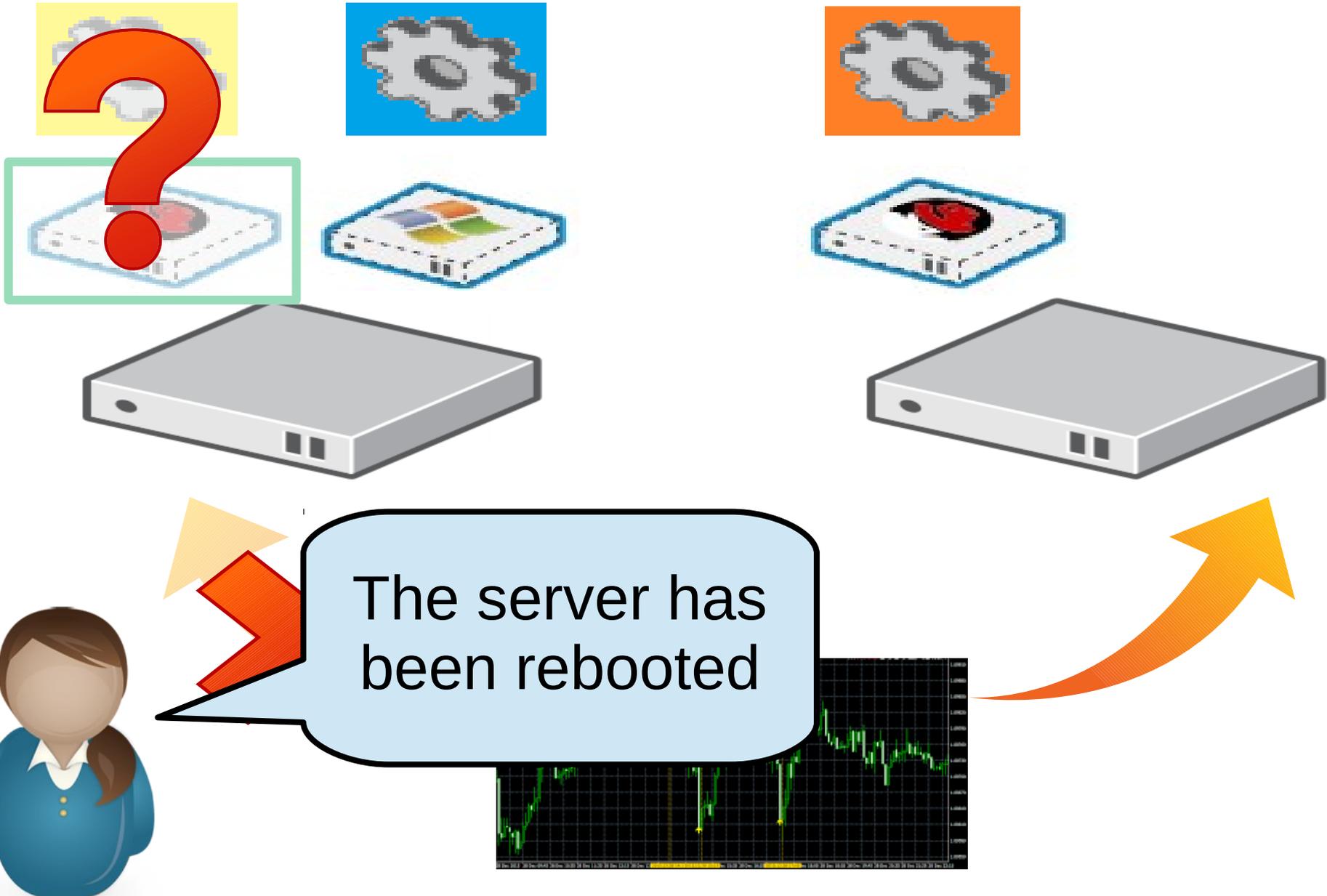
# oVirt Fault Detection – Even More Complex



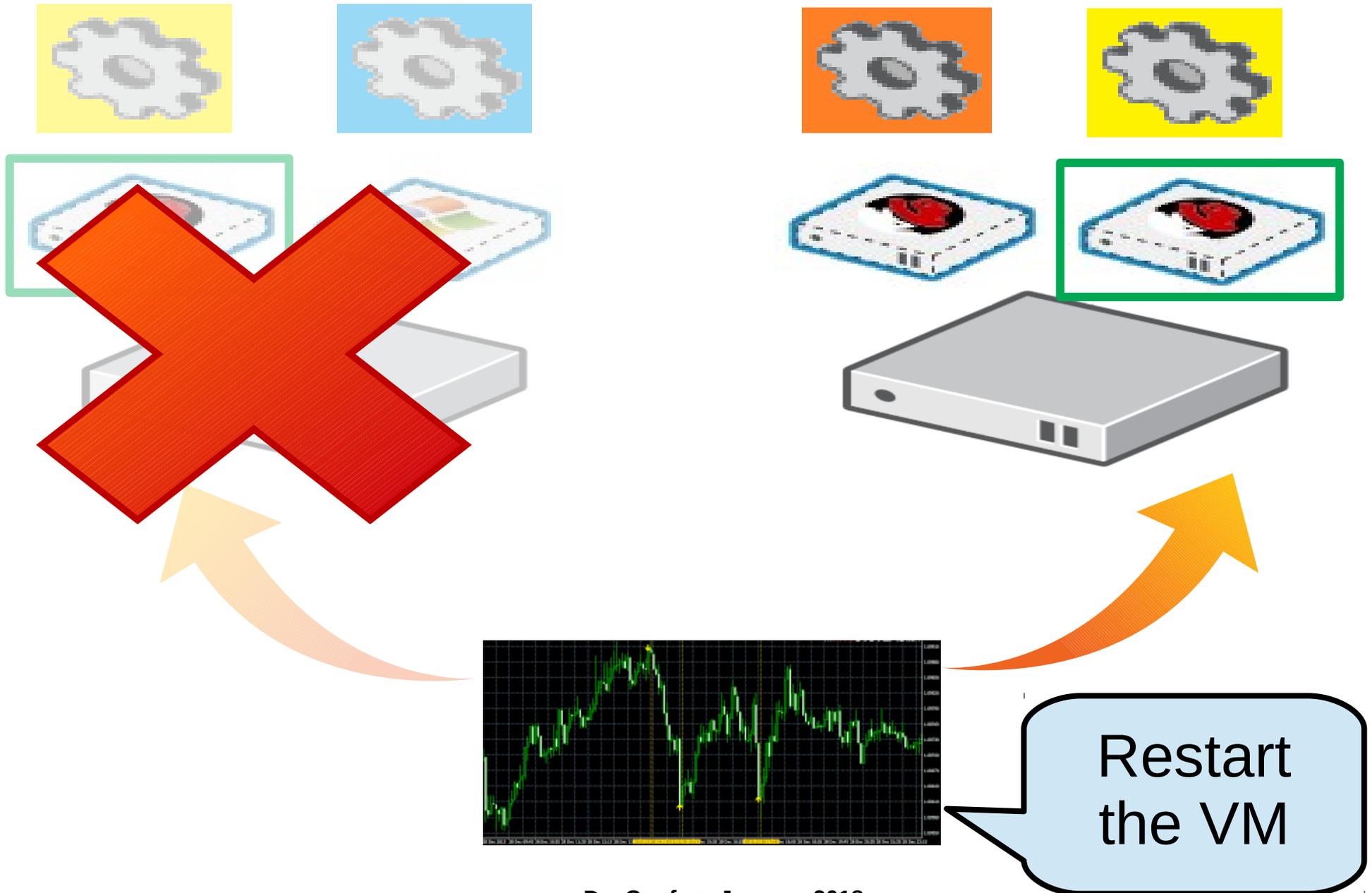
# oVirt Fault Detection – Even More Complex



# oVirt Fault Detection – Manual Confirmation

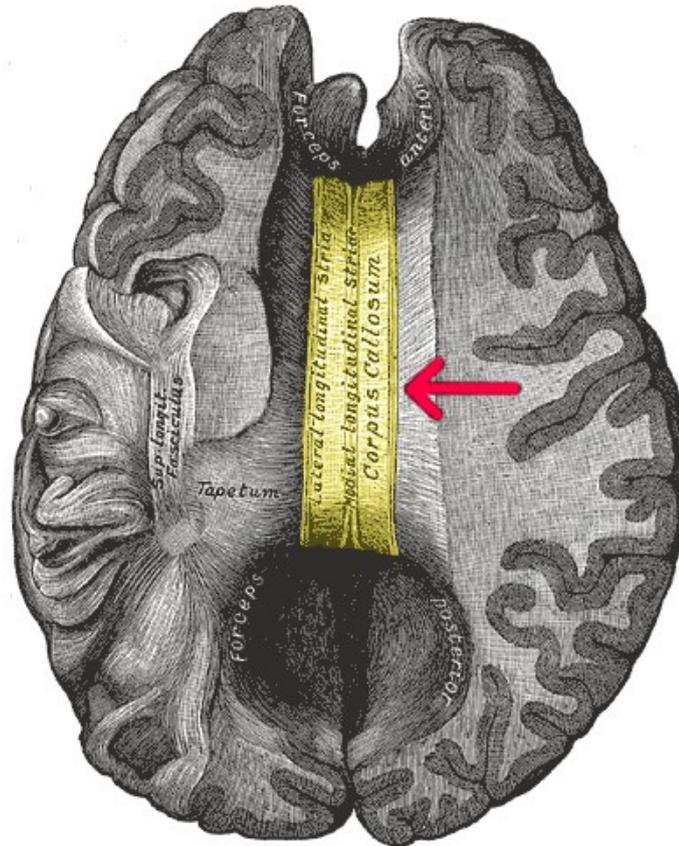


# oVirt Fault Detection – Manual Confirmation

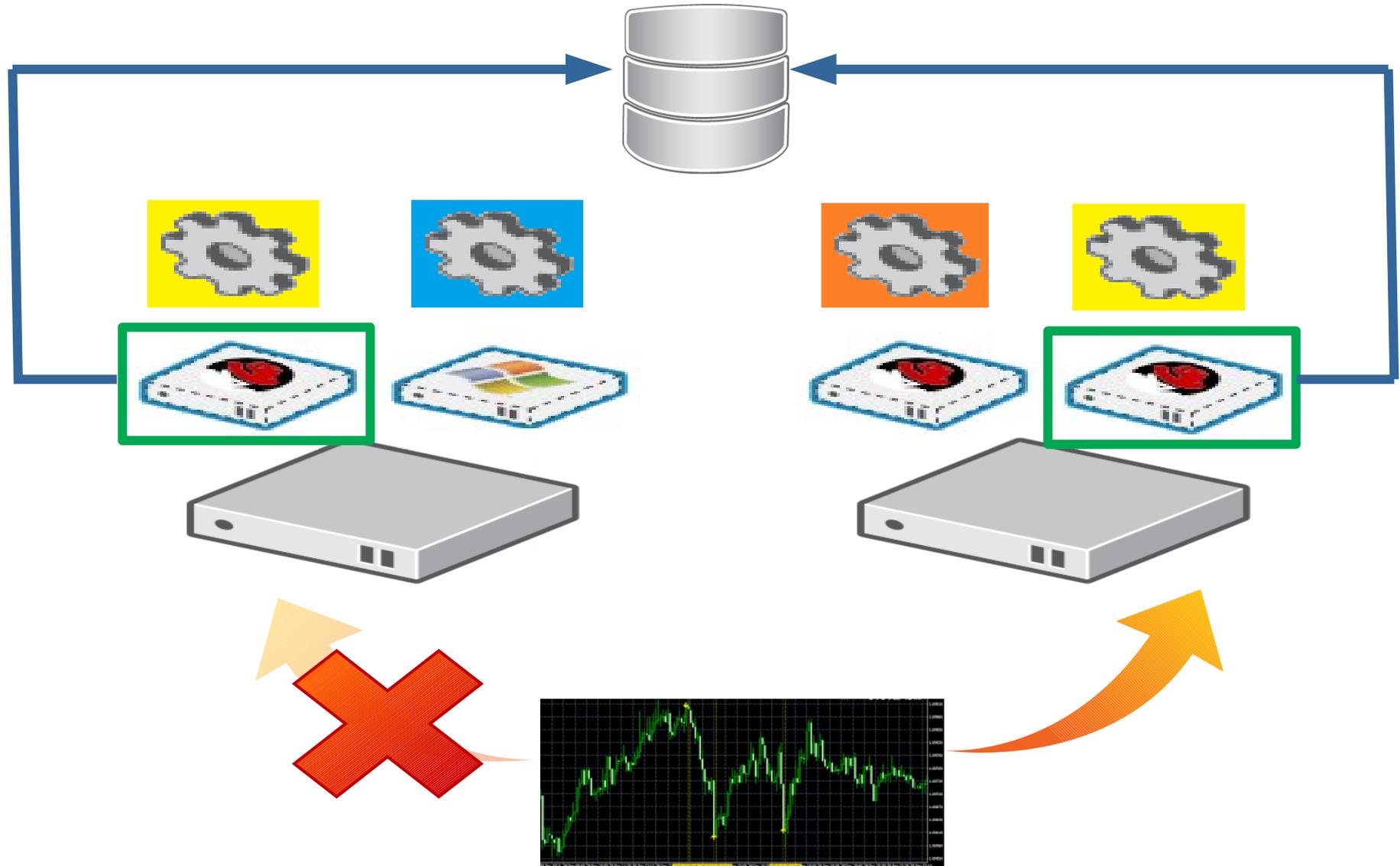


- Slow
- Error-prone
  - Mistakes may lead to a split-brain

A scenario in which several instances of the same VM run simultaneously



# Split Brain Due to a False Confirmation



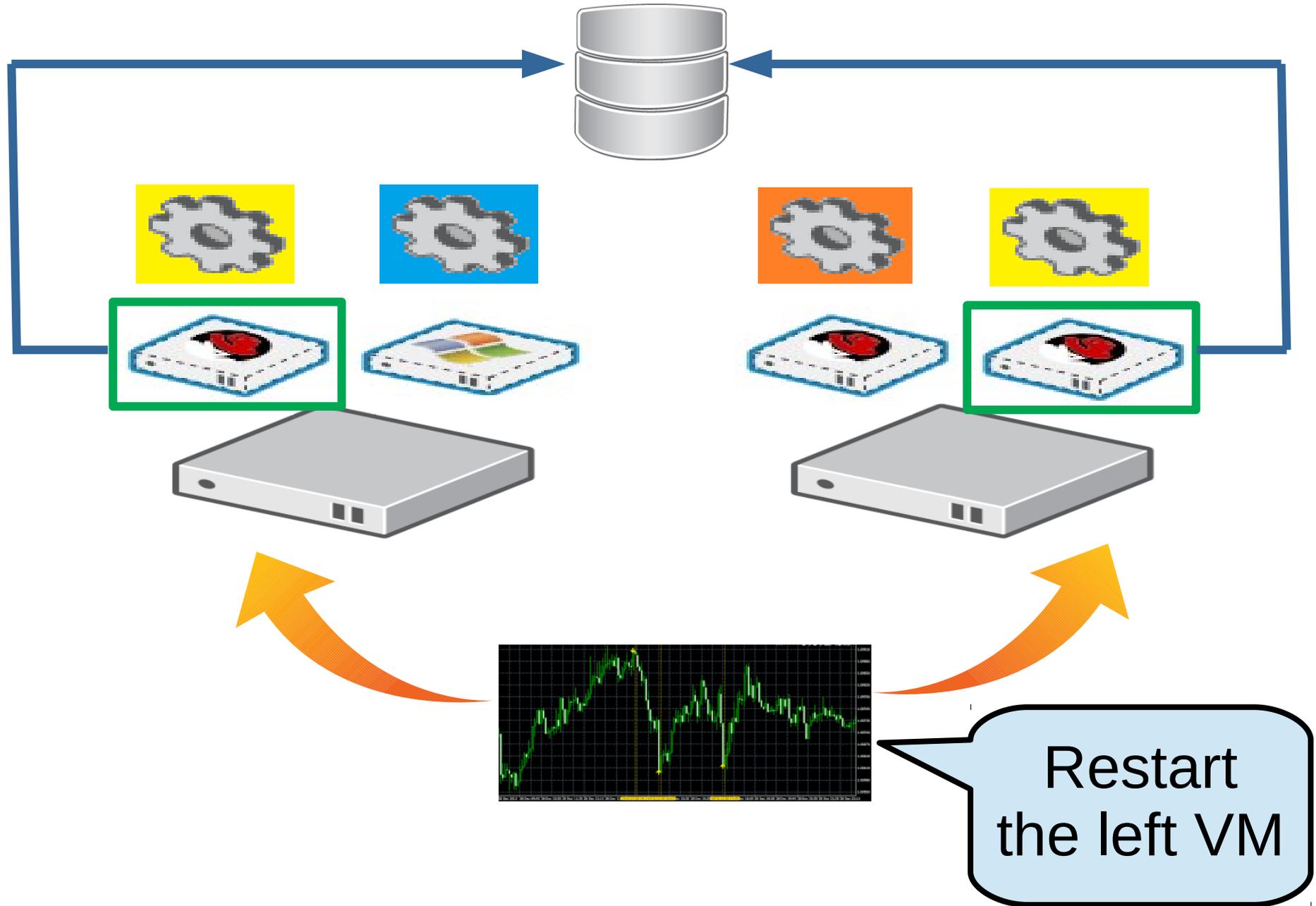
**May lead to data corruption!**

# Split Brains May Happen Due to Bugs

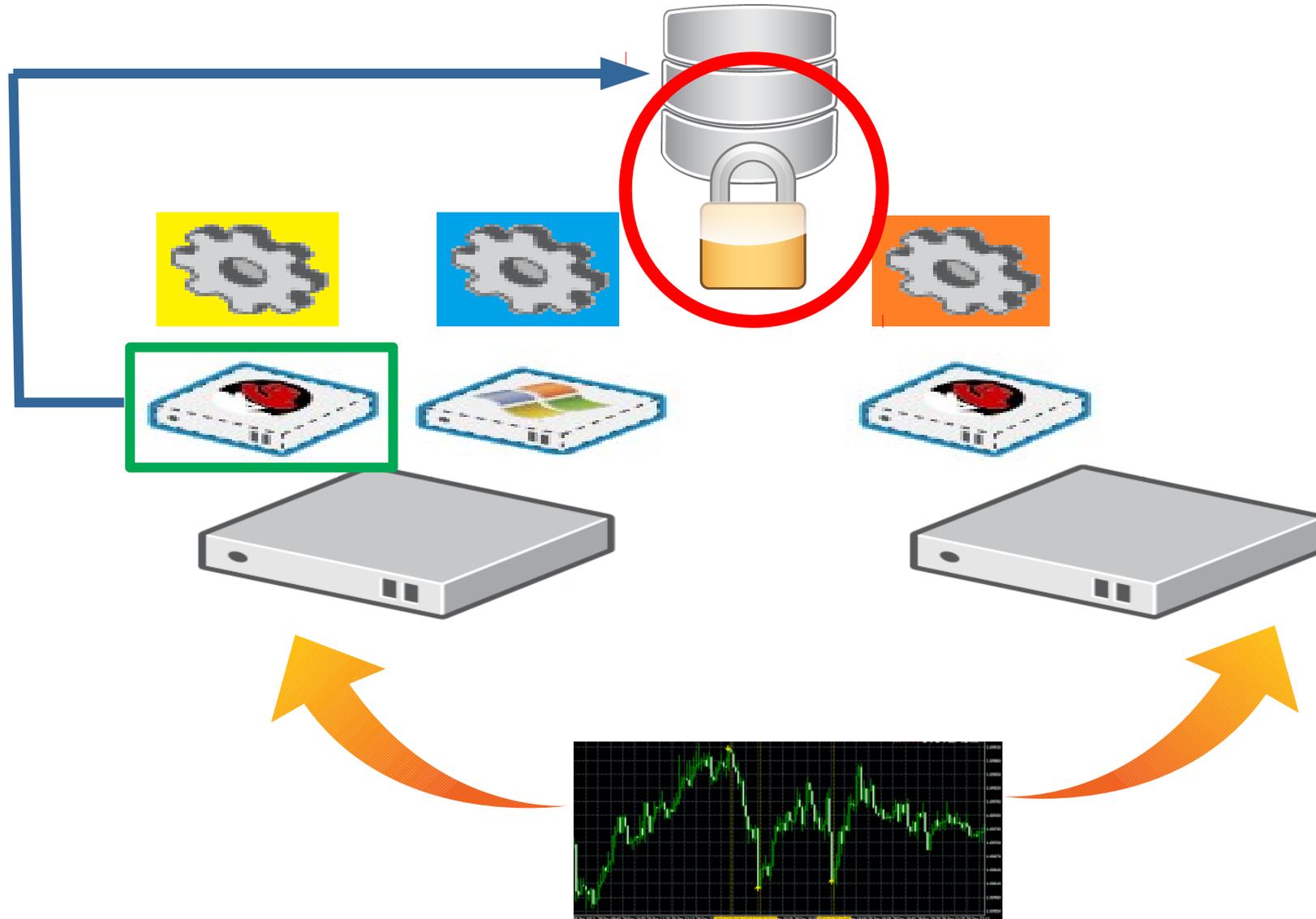


Only the right VM is reported

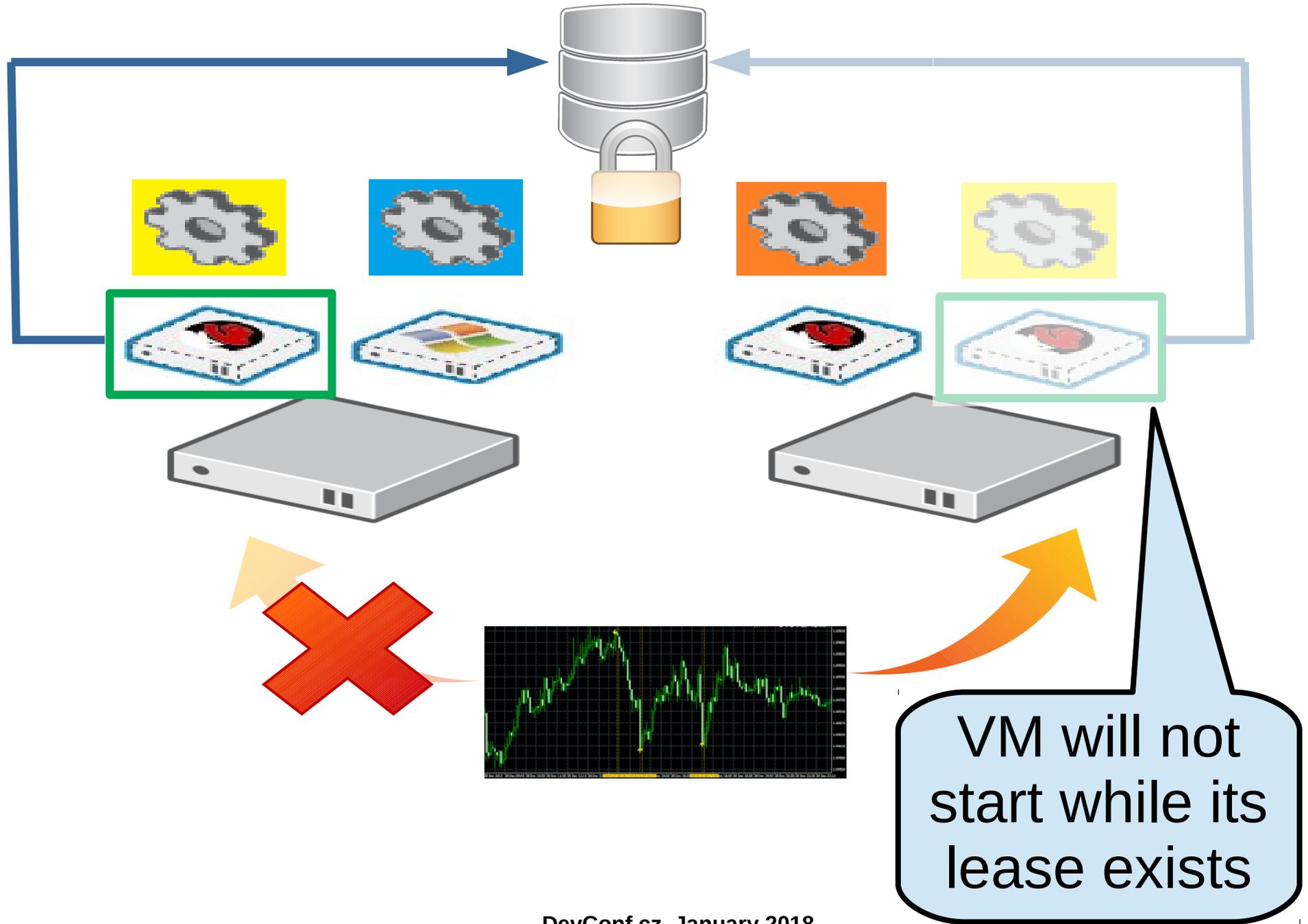
# oVirt Split Brains May Happen Due to Bugs



# oVirt VM Leases: Our Solution to Split Brains



# oVirt VM Leases: Our Solution to Split Brains



### Edit Virtual Machine ✕

<b>General</b>	Cluster	Default
<b>System</b>		<i>Data Center: Default</i>
<b>Initial Run</b>	Template	Blank   (0)
<b>Console</b>	Operating System	Debian 7
<b>Host</b>	Instance Type	Custom
	Optimized for	Server
<b>High Availability</b> >	<input checked="" type="checkbox"/> Highly Available	
<b>Resource Allocation</b>	Target Storage Domain for VM Lease	Default
<b>Boot Options</b>	Resume Behavior	KILL
<b>Random Generator</b>	<b>Priority for Run/Migration queue:</b>	
<b>Custom Properties</b>	Priority	Low
<b>Icon</b>	<b>Watchdog</b>	
<b>Foreman/Satellite</b>	Watchdog Model	No-Watchdog
<b>Affinity Labels</b>	Watchdog Action	none

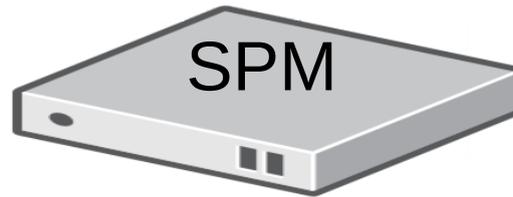
Highly Available 

Target Storage Domain for VM Lease

Default

Resume Behavior

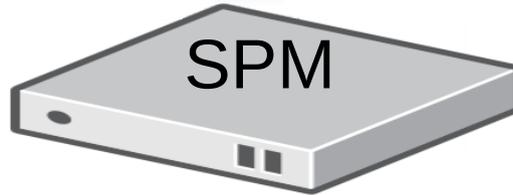
KILL



“Create a VM Lease for  
VM X in storage domain Y”



“Create a Lease X in lockspace Y”

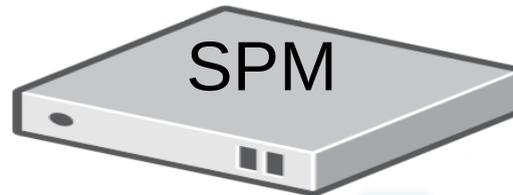


“Create a VM Lease for VM X in storage domain Y”



# VM Lease Creation

“Create a Lease X in lockspace Y”



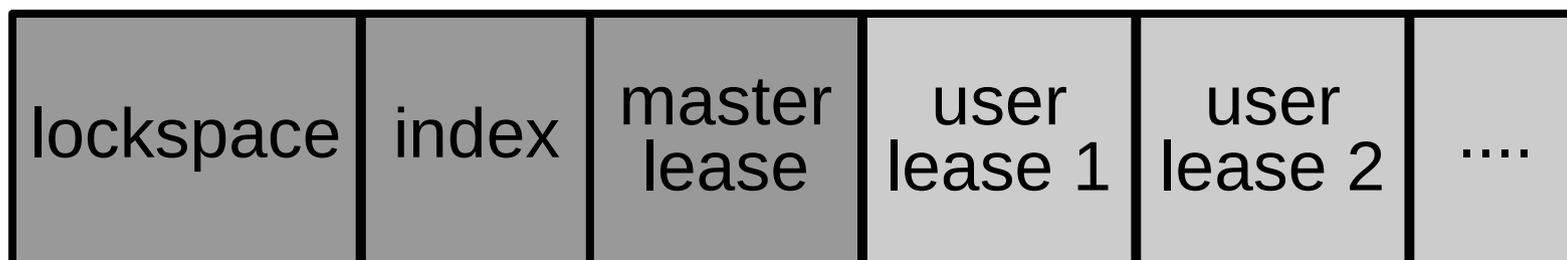
“Create a VM Lease for VM X in storage domain Y”



“Path P to xleases volume and Lease offset O”



- Sanlock does not manage leases allocation
- Volume layout:



- Same format in block and file storage
- [Deep Dive - VM leases](#) (youtube)

```
<domain type='kvm' id='6'>
```

```
<name>fedora8</name>
```

```
... skipped ...
```

```
<devices>
```

```
... skipped ...
```

```
<lease>
```

```
<lockspace>571184ae-79da-41fb-a3fb-c3117991abae</lockspace>
```

```
<key>cbd783e4-45f8-4b51-93ca-4460d4dad772</key>
```

```
<target path='/rhev/data-center/mnt/10.35.1.90:_srv_Default/571184ae-79da-41fb-a3fb-c3117991abae/dom_md/xleases' offset='3145728' />
```

```
</lease>
```

```
... skipped ...
```

```
</domain>
```

# oVirt Running a VM with a Lease

oVirt

Acquires the Lease using Sanlock



Domain XML with Lease

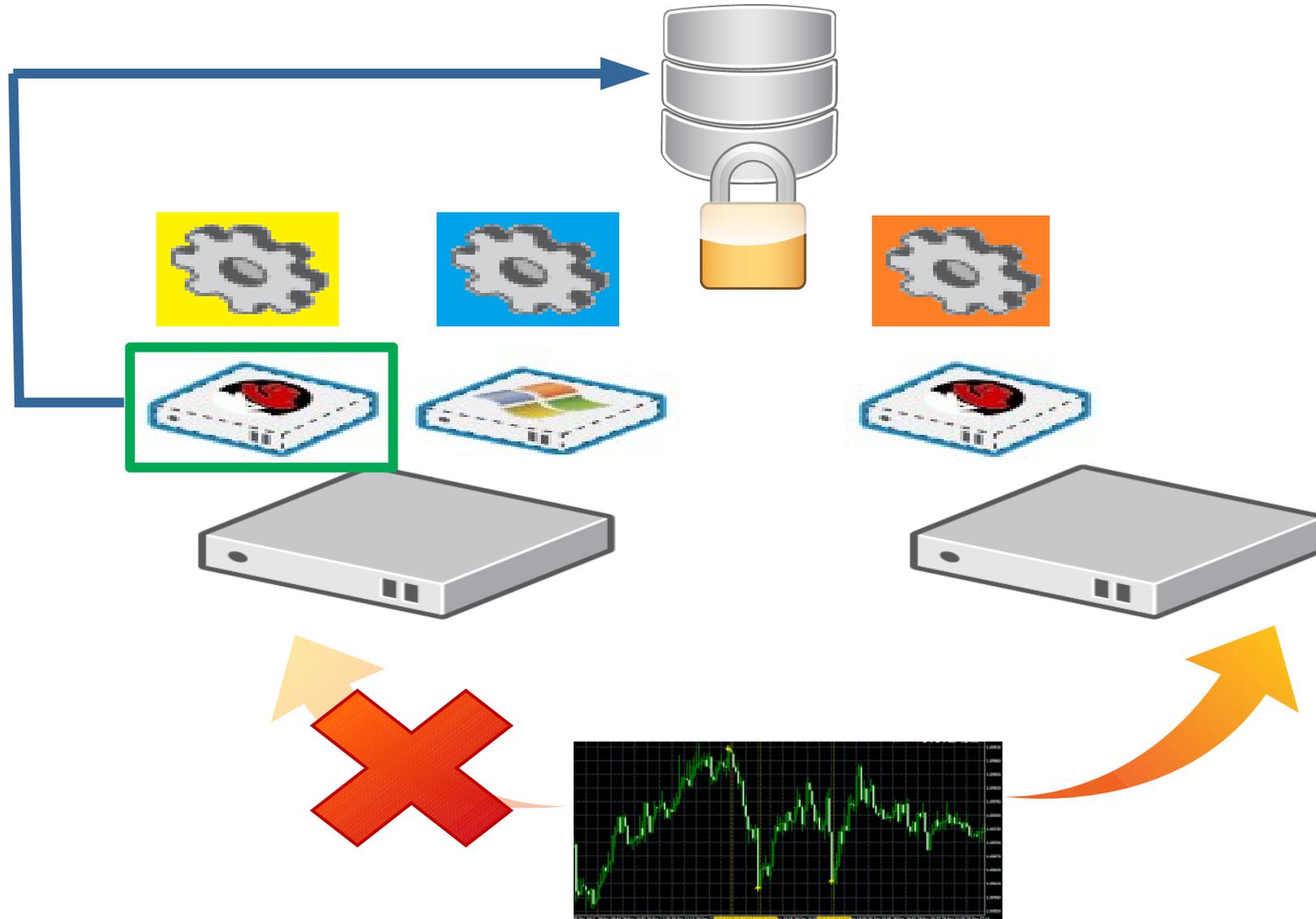
Lease



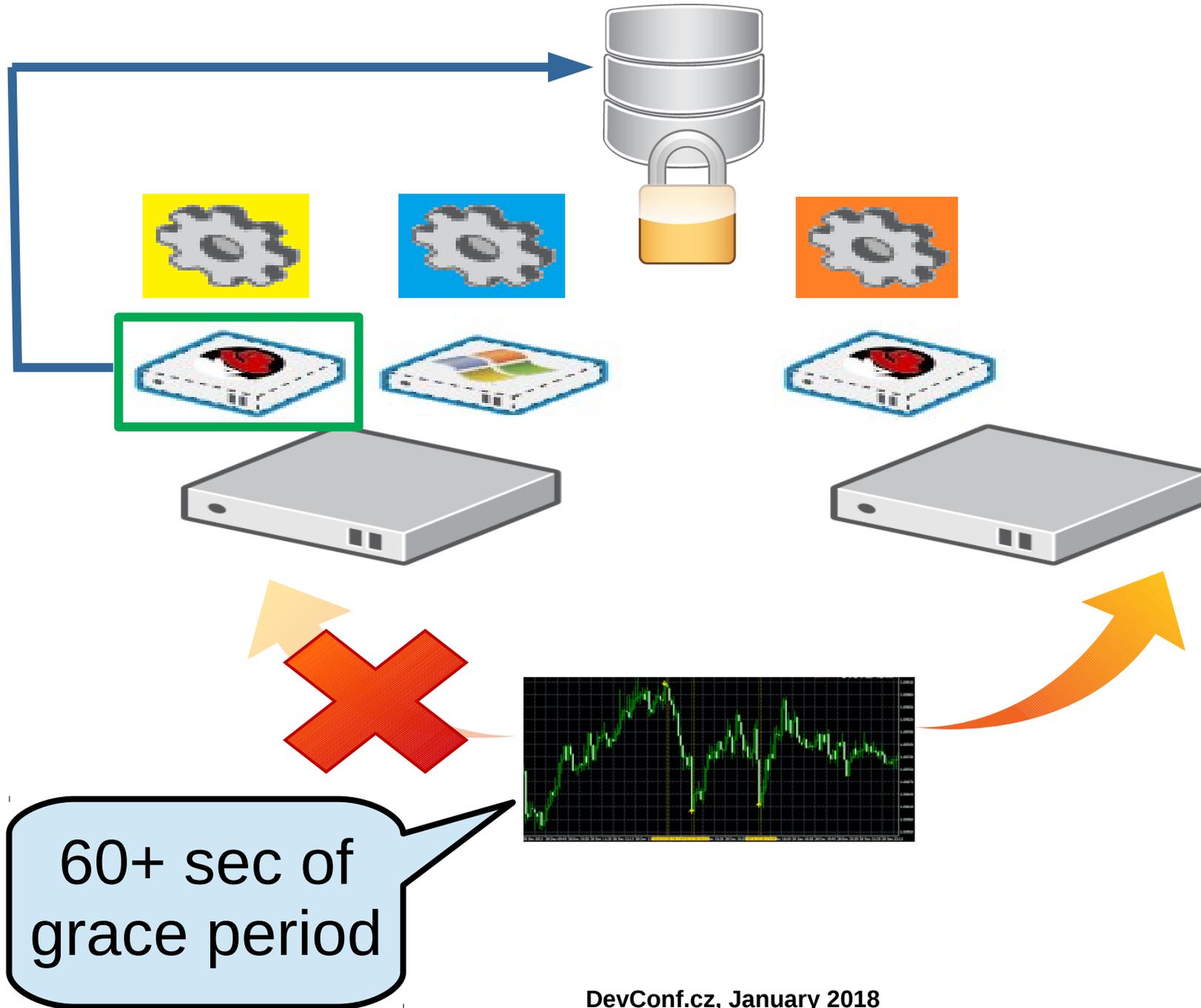
# oVirt Non-Responsive Host Treatment



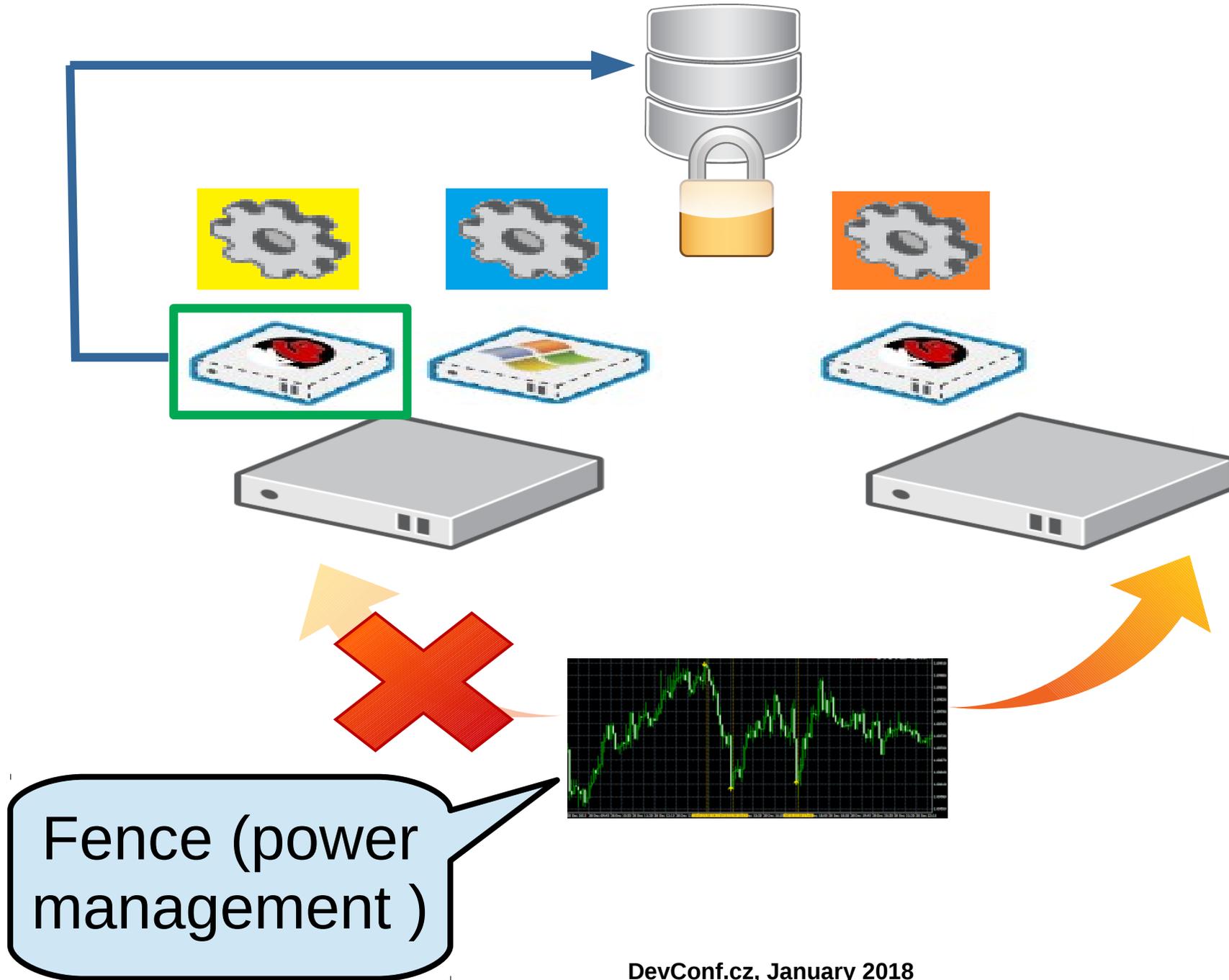
# oVirt Non-Responsive Host Treatment



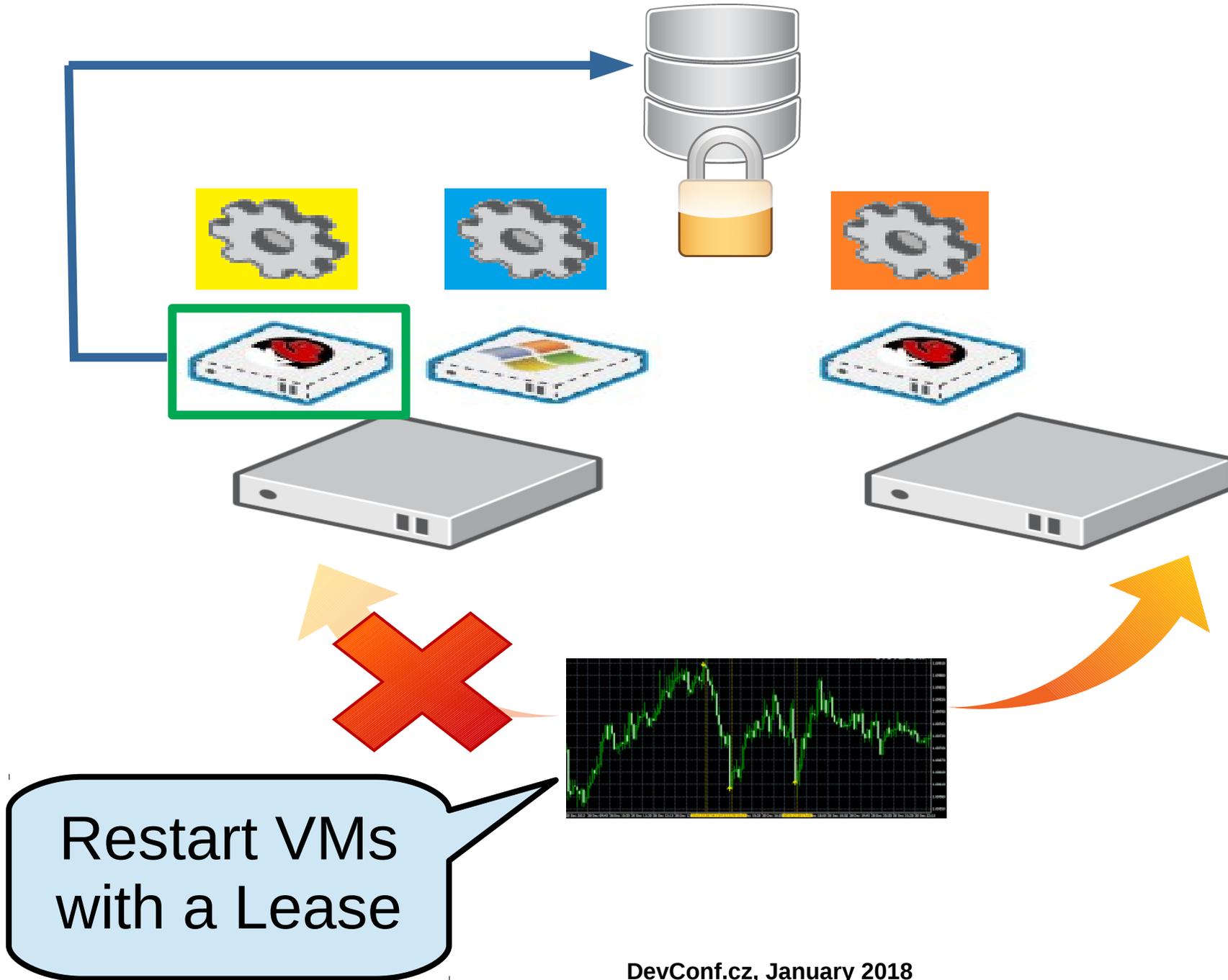
# oVirt Non-Responsive Host Treatment



# oVirt Non-Responsive Host Treatment



# oVirt Non-Responsive Host Treatment

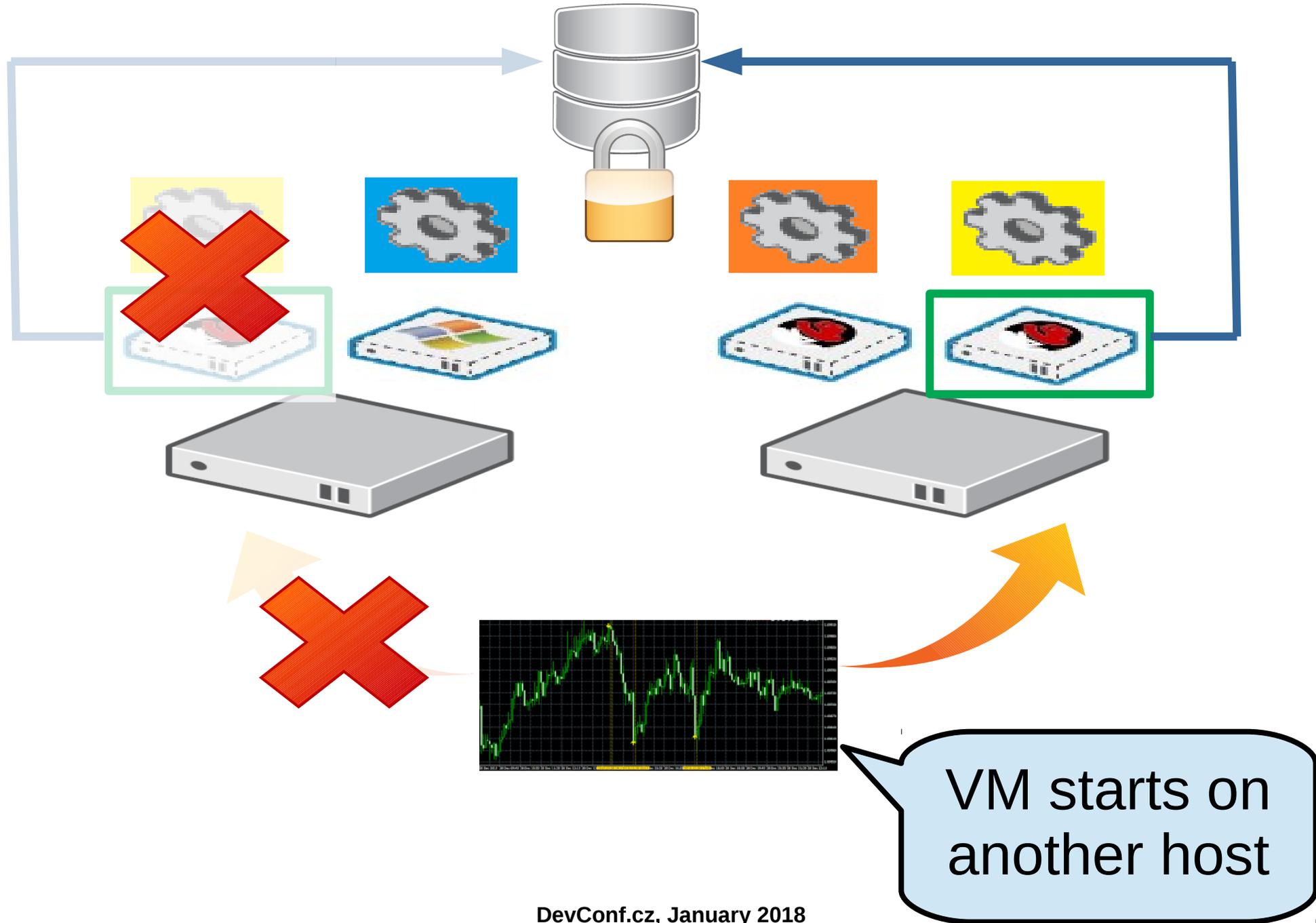


# (1) Non-Responsive Host + VM is Down

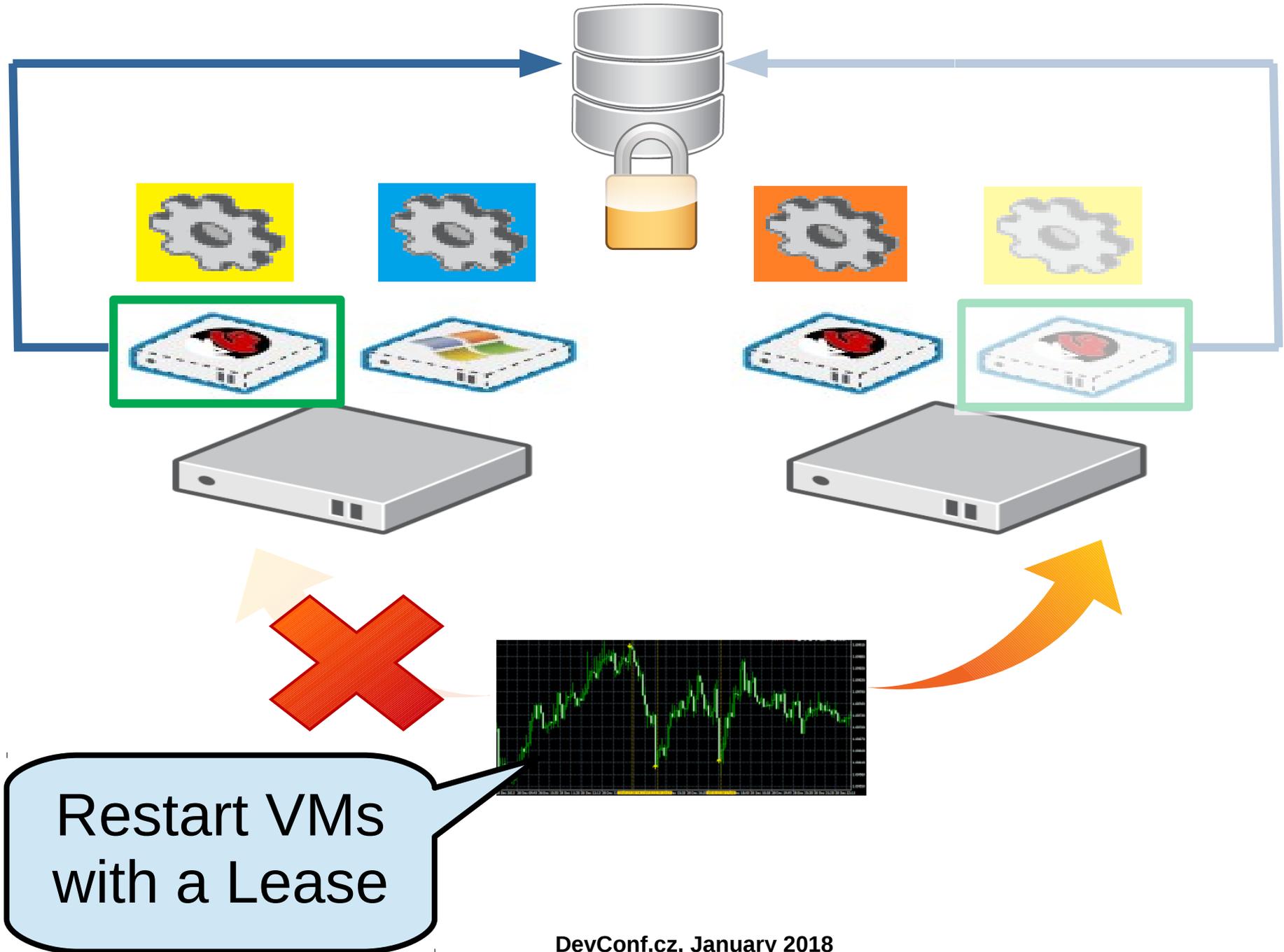


Restart VMs with a Lease

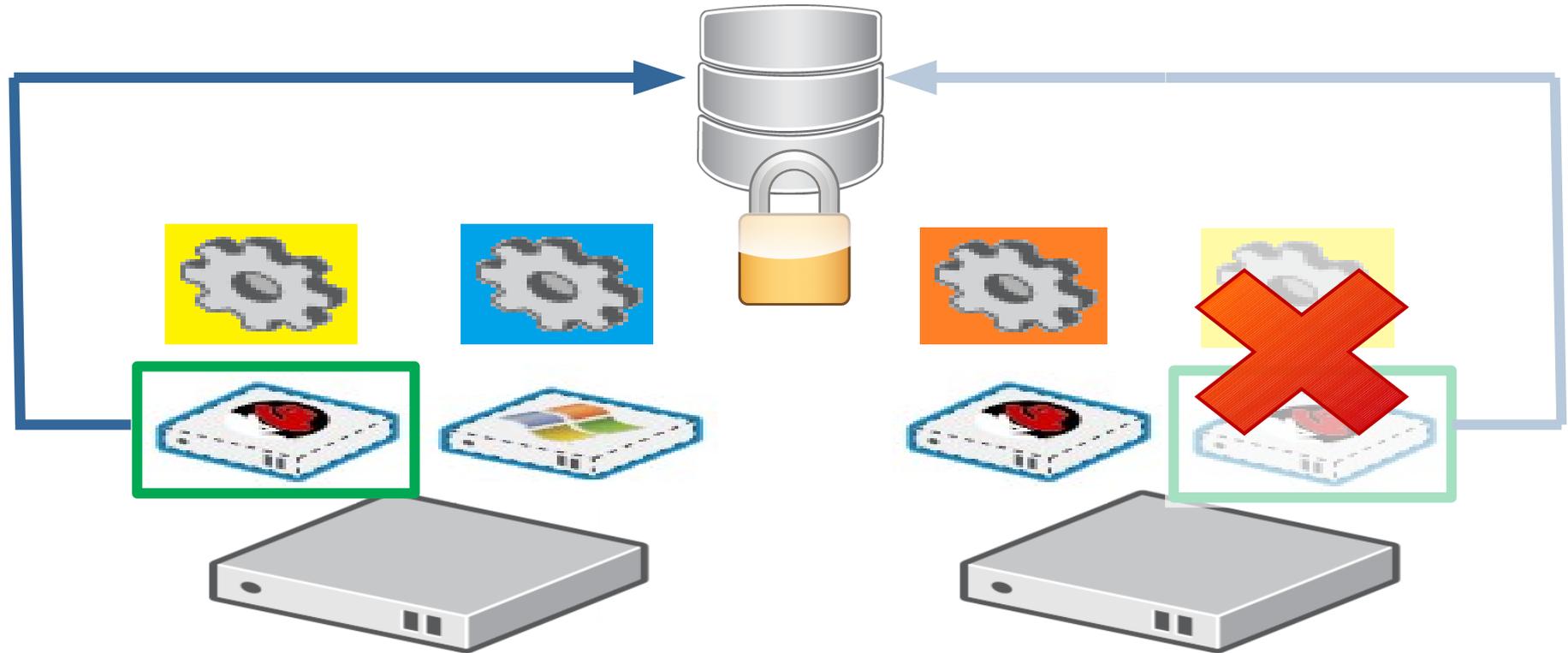
# (1) Non-Responsive Host + VM is Down



# (2) Non-Responsive Host + VM is UP



# (2) Non-Responsive Host + VM is UP

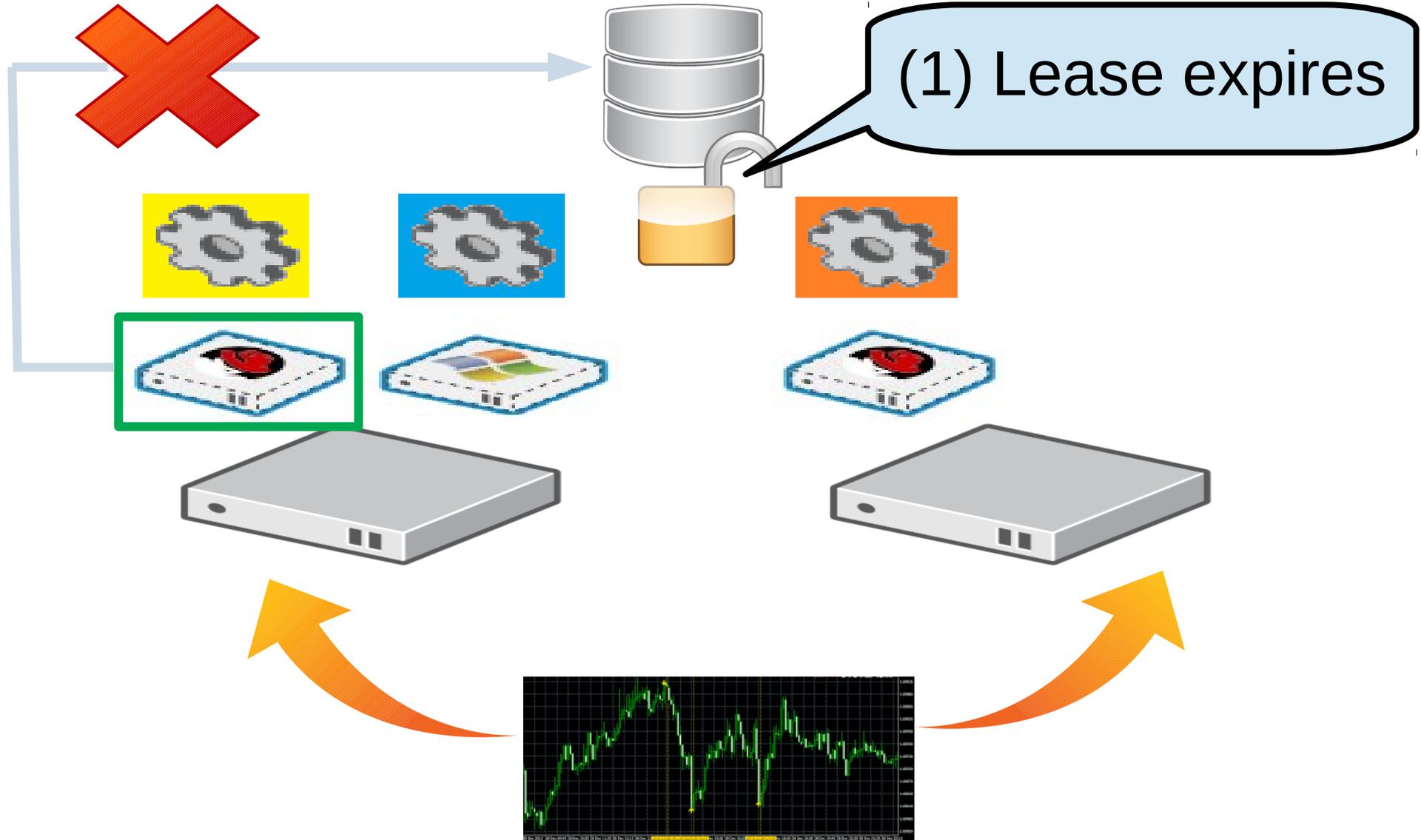


Restart VMs  
with a Lease

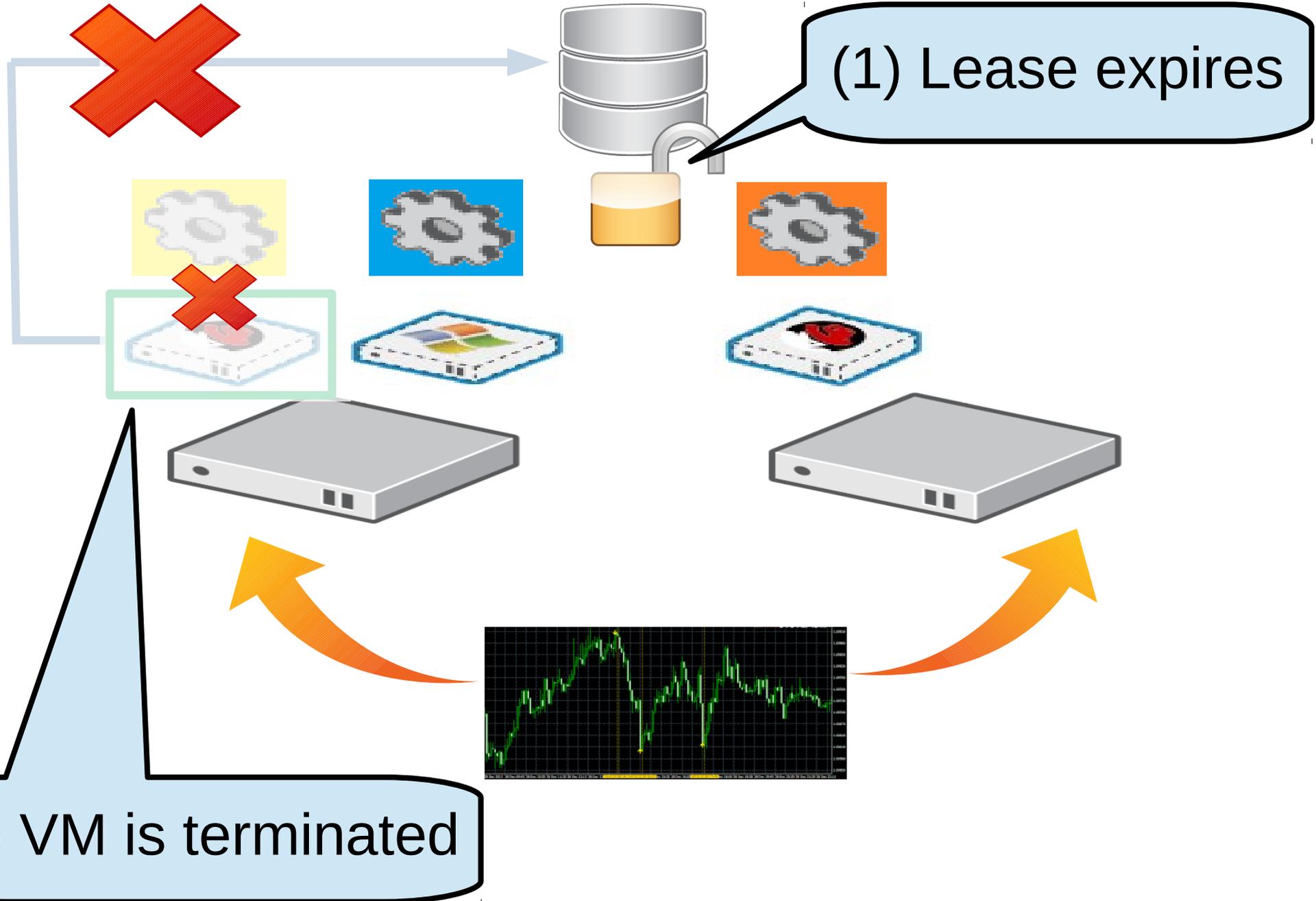
# Disconnection From Storage Device



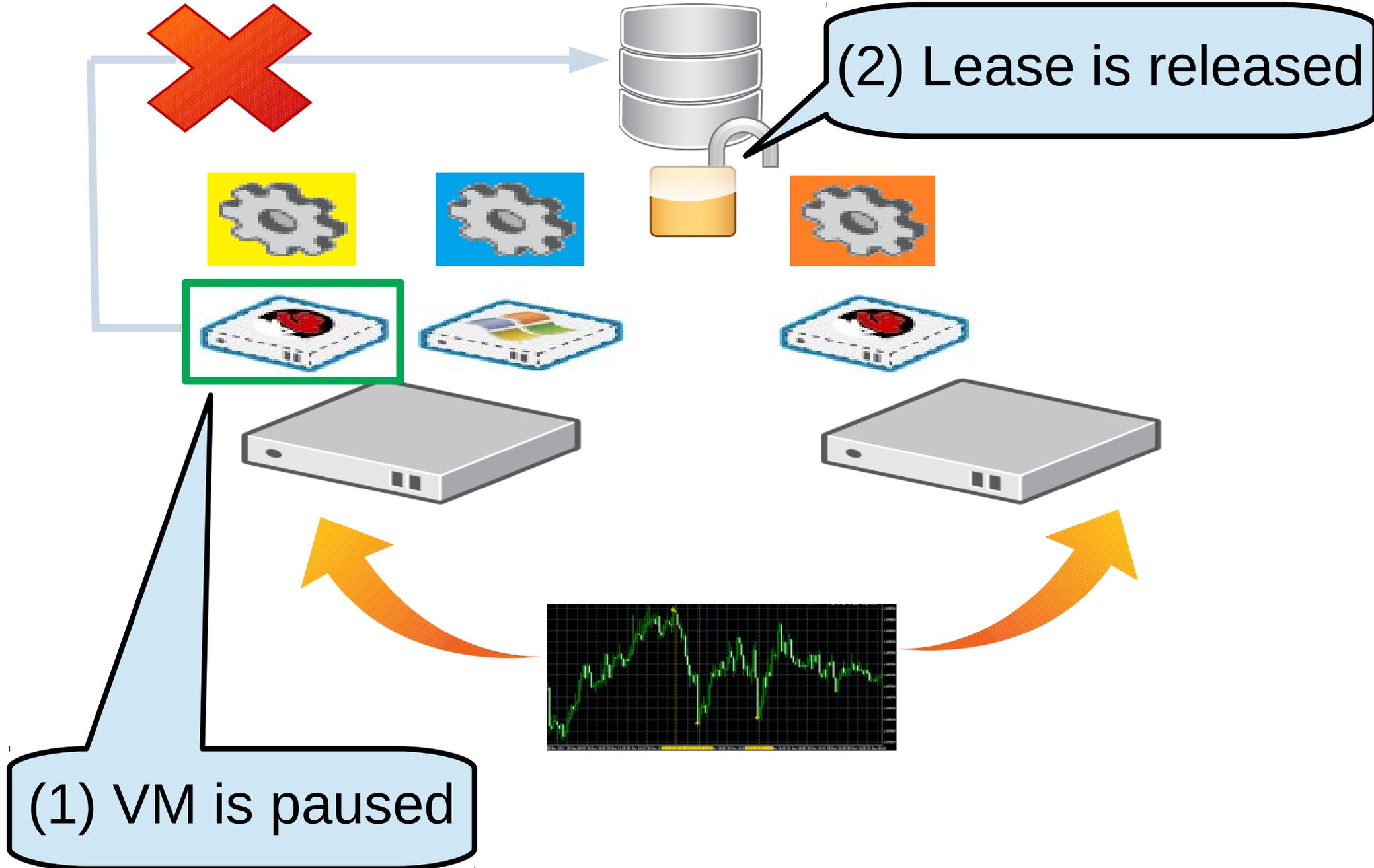
# Disconnection From Storage Device (1)



# Disconnection From Storage Device (1)



# Disconnection From Storage Device (2)



- VM Lease – an important new element
  - Prevents split-brains
  - Enables automatic restart of unreported VMs
- Available since oVirt 4.1
  - Polished in oVirt 4.2
- Possible future enhancements:
  - May be used to restart paused VMs
  - Move together with the bootable disk

# THANK YOU!

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