# High Availability with
# No Split Brains!
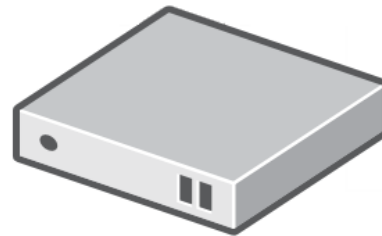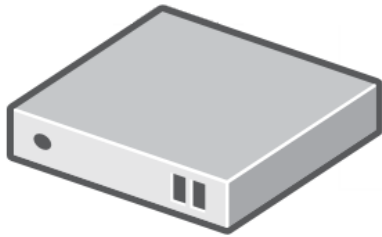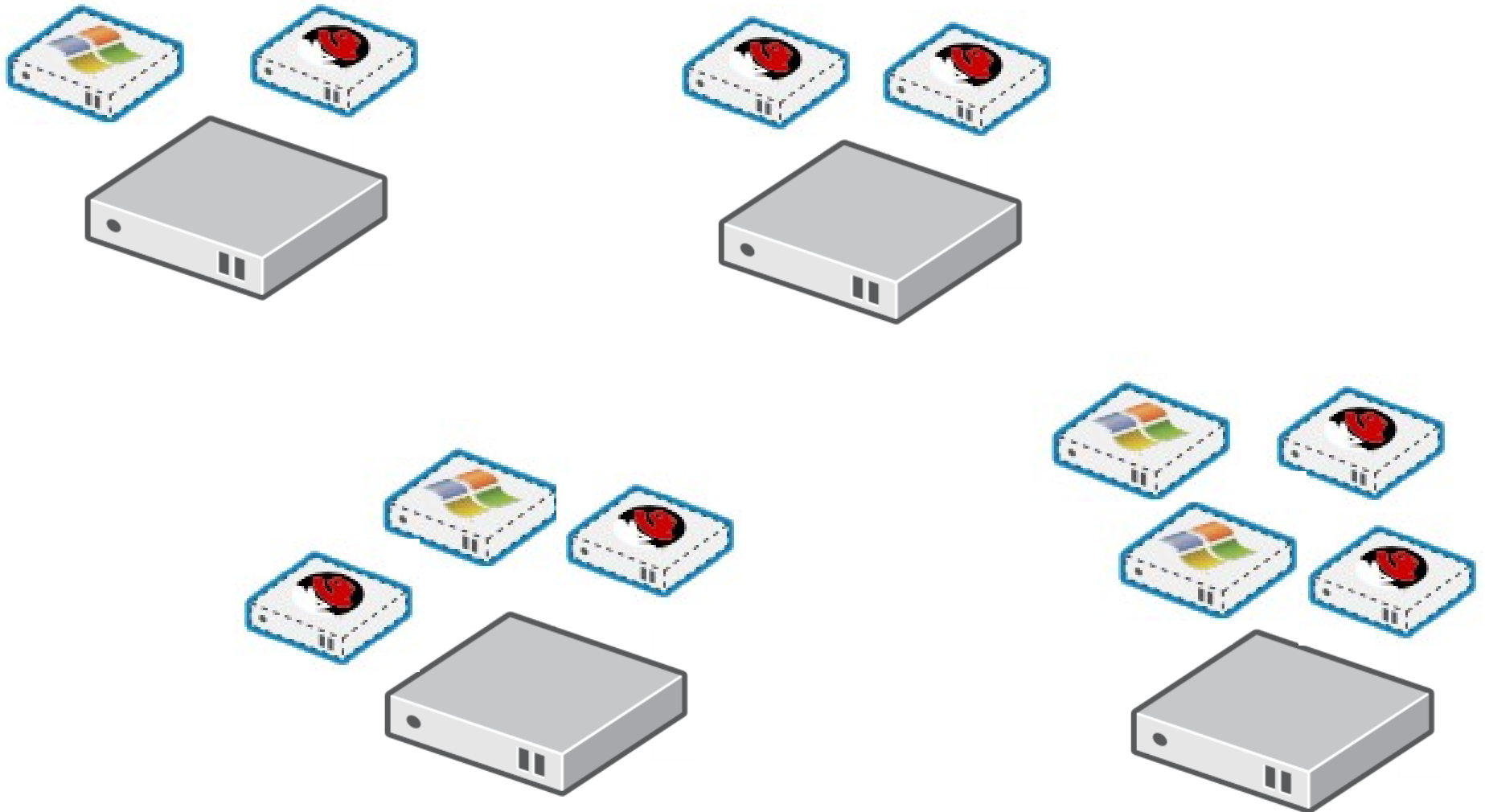
Arik Hadas
Principal Software Engineer
Red Hat
27/01/2018

- Higher resource consumption

- More responsibility on the application

- Backup lives in a different environment

  - Different IP address(es)

  - Different disk(s)

- More efficient resource consumption

- Implemented at the infrastructure level

- Backup starts in the same environment
  - Same IP address(es)
  - Same disk(s)

HA VM went down!

What if:

- Inaccessible resources

- VM is locked

- VM is being
  intentionally shut down

Restart
the VM

# Automatic Restart – Not That Simple

What if:
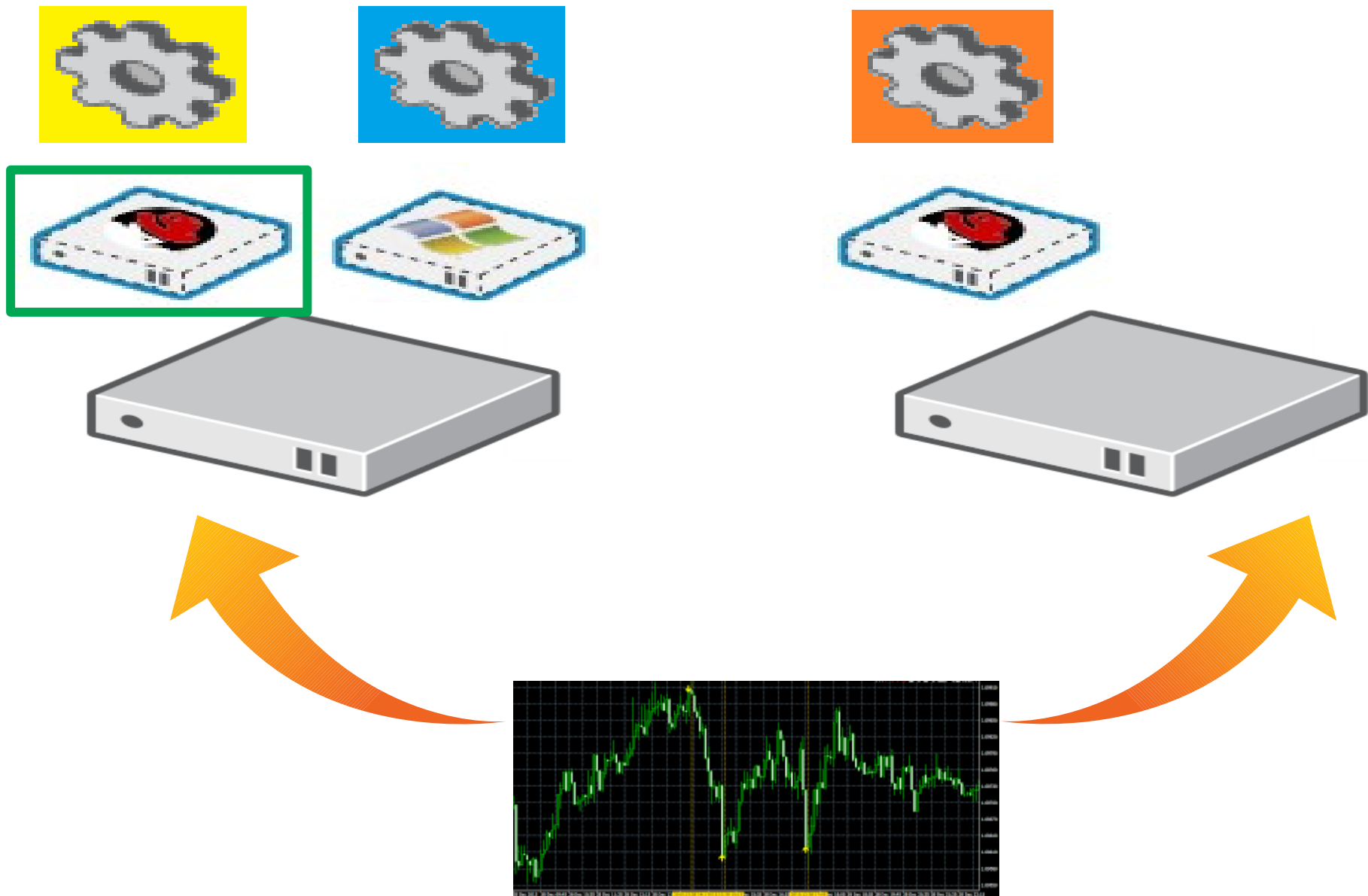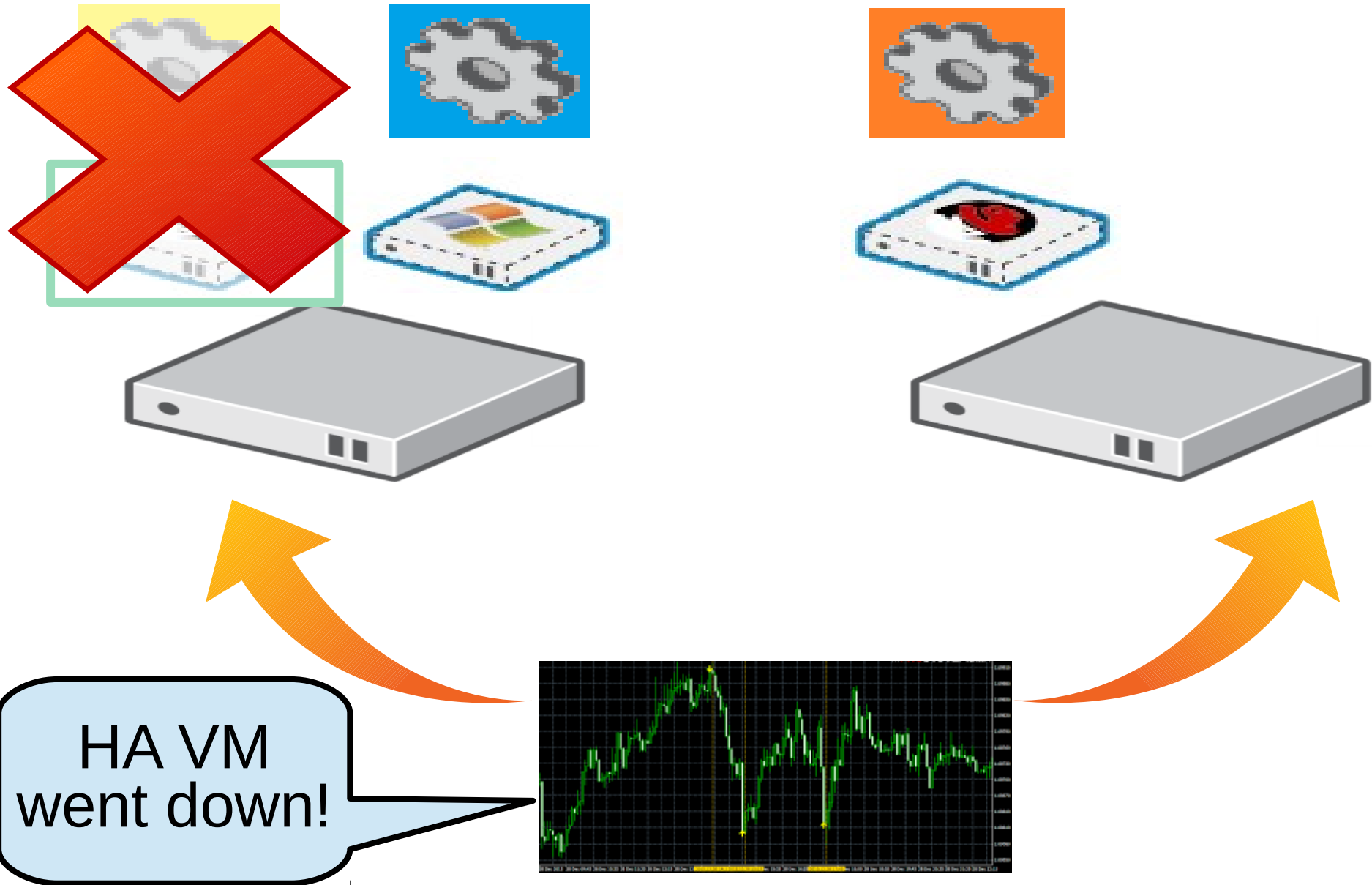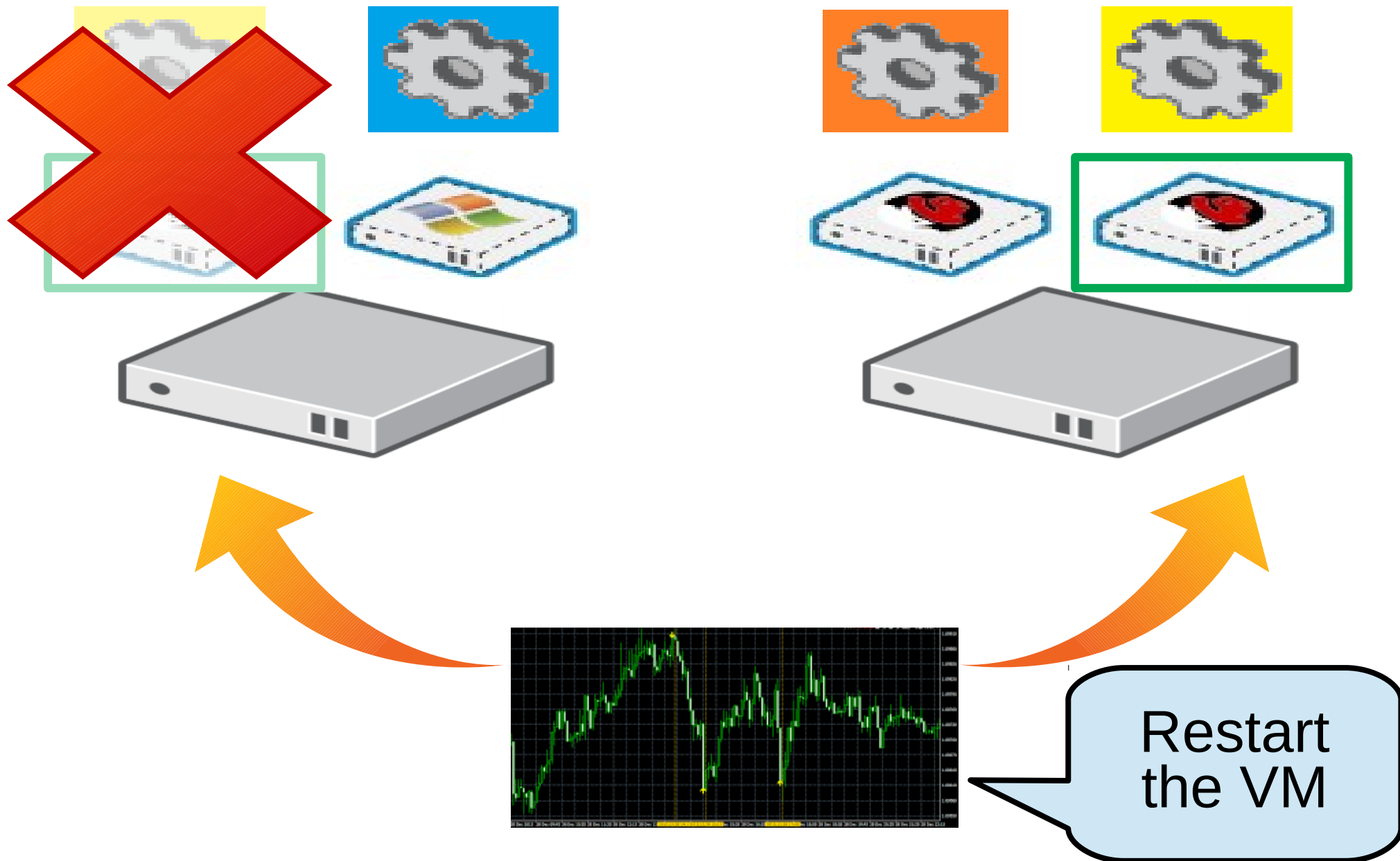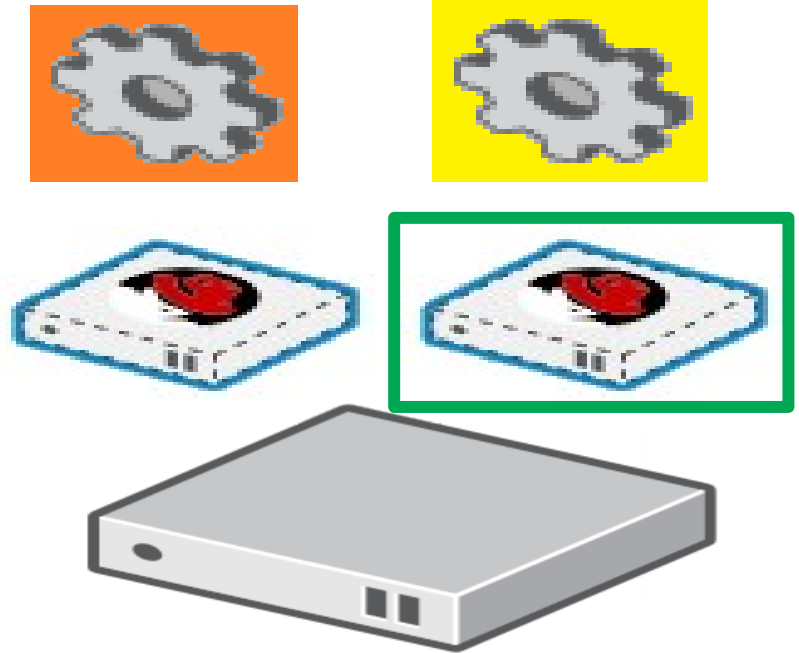
– Inaccessible resources

– 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

The server has been rebooted

Restart
the VM

- Slow

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

A scenario in which several VMs that may write to the same disk(s) are running simultaneously

# May lead to data corruption!

Only the right VM is reported (on startup)

# Split Brains May Happen Due to Bugs

Restart the left VM

DevConf.cz, January 2018

VM will not start while its lease exists

# VM Lease Creation

☑ Highly Available

Target Storage Domain for VM Lease

Default ⌄

Resume Behavior

KILL ⌄

SPM

"Create VM Lease for
 VM X in storage domain Y"

**DevConf.cz, January 2018**

"Create a Lease X in
lockspace Y"

SPM

"Create VM Lease for
VM X in storage domain Y"

oVirt

"Create a Lease X in lockspace Y"

SPM

"Create VM Lease for VM X in storage domain Y"

"Path P to xleases volume and Lease offset O"

# xleases volume

- Sanlock does not manage leases allocation

- Volume layout:

| lockspace | index | master lease | user lease 1 | user lease 2 | .... |
|-----------|-------|--------------|--------------|--------------|------|

- 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>
```
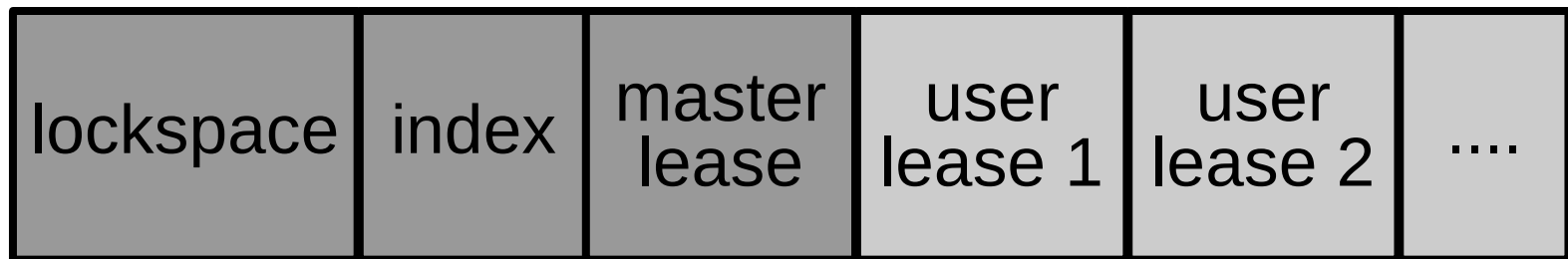
60+ sec of grace period

Fence (power Management )

Restart VMs with a Lease

# Non-Responsive Host + VM is Down



VM starts on another host

oVirt

Restart VMs with a Lease

Restart VMs with a Lease

(1) Lease expires

(2) VM is terminated

(2) Lease is released

(1) VM is paused

# Summary

- 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

- Future enahncements:

    – 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