# End of the Road - Facing Current Scaling Limits within OpenStack

OpenStack Summit, Vancouver                                    May 19, 2015

Christian Berendt
Cloud Solution Architect
B1 Systems GmbH
berendt@b1-systems.de

Thomas Kaergel
Linux Consultant & Developer
B1 Systems GmbH
kaergel@b1-systems.de

## Introducing B1 Systems

- founded in 2004
- operating both nationally and internationally
- more than 60 employees; low employee turnover
- Provider for IBM, SUSE, Oracle & HP
- vendor-independent (hardware and software)
- Focus:
    - Consulting
    - Support
    - Development
    - Training
    - Operations
    - Solutions

# Areas of Expertise

- Virtualization (XEN, KVM & RHEV)
- Systems management (Spacewalk, Red Hat Satellite, SUSE Manager)
- Configuration management (Puppet & Chef)
- Monitoring (Nagios & Icinga)
- IaaS Cloud (OpenStack & SUSE Cloud)
- High availability (Pacemaker)
- Shared Storage (GPFS, OCFS2, DRBD & CEPH)
- File Sharing (ownCloud)
- Packaging (Open Build Service)
- Providing on-site systems administration and/or development

Once upon a time ...

Source: lassedesignen/Shutterstock.com

Source: Ken Pepple, http://ken.pepple.info/

**Source: varunsingh180000/Pixabay.com**

# What Happenend?

## Observations

- nova list extremely slow
- almost all nova operations on instances affected
- horizon too slow to be usable
- DB and nova services under heavy load

# Observations

- nova list extremely slow
- almost all nova operations on instances affected
- horizon too slow to be usable
- DB and nova services under heavy load

## Observations

- nova list extremely slow
- almost all nova operations on instances affected
- horizon too slow to be usable
- DB and nova services under heavy load

# Observations

- nova list extremely slow
- almost all nova operations on instances affected
- horizon too slow to be usable
- DB and nova services under heavy load
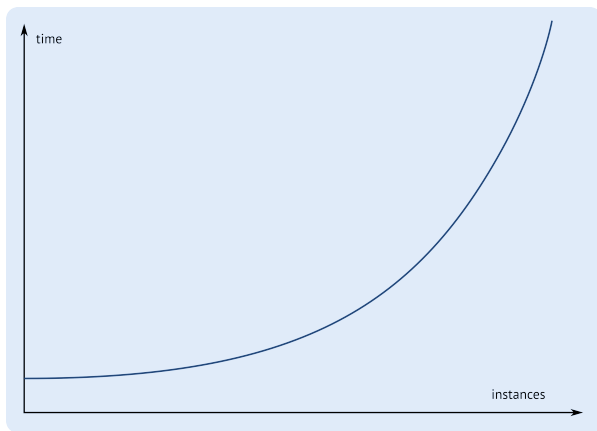
Case Study

# Many Instances in Single Tenant (Folsom)



Figure : nova-list Duration over Instance Count
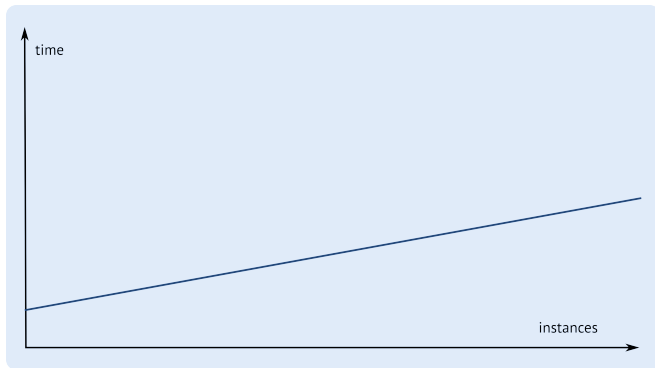
# Many Instances in Multiple Tenants (Folsom)



Figure : nova-list Duration over Instance Count
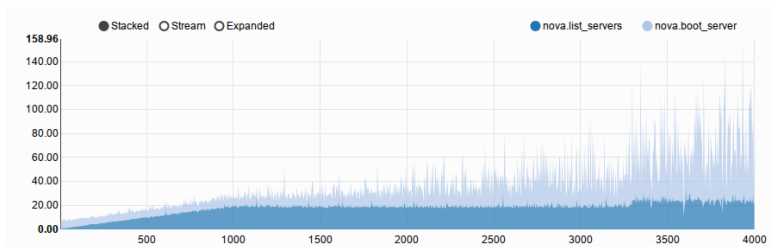
# Many Instances in Single Tenant (Today)



Figure : nova-list Duration over Instance Count

# Investigation Strategy



Source: OpenClips/Pixabay.com

## Actions

- watch CPU load on infrastructure during load situation
- switch logmode to debug
- observe logs during load situation
- turn mysql query logging on and watch the DB queries
- analyze the code

## Actions

- watch CPU load on infrastructure during load situation
- switch logmode to debug
- observe logs during load situation
- turn mysql query logging on and watch the DB queries
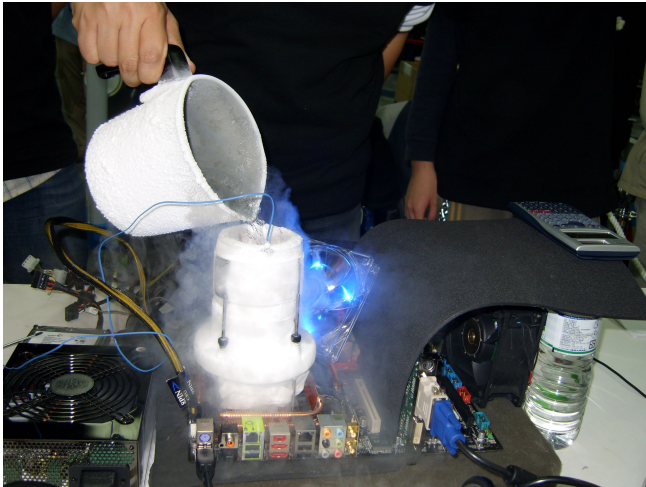- analyze the code

## Actions

- watch CPU load on infrastructure during load situation
- switch logmode to debug
- observe logs during load situation
- turn mysql query logging on and watch the DB queries
- analyze the code

## Actions

- watch CPU load on infrastructure during load situation
- switch logmode to debug
- observe logs during load situation
- turn mysql query logging on and watch the DB queries
- analyze the code

## Actions

- watch CPU load on infrastructure during load situation
- switch logmode to debug
- observe logs during load situation
- turn mysql query logging on and watch the DB queries
- analyze the code

## sqlalchemy DB-Join, linewidth 80 characters, length 100+ lines

```
SELECT instances.created_at AS instances_created_at,
instances.updated_at AS instances_updated_at,
instances.deleted_at AS instances_deleted_at,
instances.id AS instances_id,
instances.user_id AS instances_user_id,
instances.project_id AS instances_project_id,
instances.image_ref AS instances_image_ref,
instances.kernel_id AS instances_kernel_id,
[...]
...
...
[...]
WHERE instances.deleted = 0 AND instances.host = 'computexen0158'
```

# Possible Solutions

Source: RicoShen/Wikimedia.org

- more powerful hardware for Nova and DB
- rewrite nova/sqlalchemy code that generates those big DB joins
- reorganize user/tenant-layout for the use case

- more powerful hardware for Nova and DB
- rewrite nova/sqlalchemy code that generates those big DB joins
- reorganize user/tenant-layout for the use case

- more powerful hardware for Nova and DB
- rewrite nova/sqlalchemy code that generates those big DB joins
- reorganize user/tenant-layout for the use case

# Prevention Strategy

- determine expected load and expected elasticity
- design for horizontal scalability using active/active HA setups whenever possible
- built a representative miniature of your cloud for measurements, experiments and development

# Useful Tools

# Vagrant



Source: Fco.plj/de.wikipedia.org

# Vagrant Advantages

- reproducible and portable work environments
- easy to set up and learn
- usable for scale testing and development
- many *providers* available (Virtualbox, KVM, VMware...) to virtualize hosts
- choice between many *provisioners* (Shell, Ansible, Chef, Puppet...) to configure hosts

# Vagrant Advantages

- reproducible and portable work environments
- easy to set up and learn
- usable for scale testing and development
- many *providers* available (Virtualbox, KVM, VMware...) to virtualize hosts
- choice between many *provisioners* (Shell, Ansible, Chef, Puppet...) to configure hosts

# Vagrant Advantages

- reproducible and portable work environments
- easy to set up and learn
- usable for scale testing and development
- many *providers* available (Virtualbox, KVM, VMware...) to virtualize hosts
- choice between many *provisioners* (Shell, Ansible, Chef, Puppet...) to configure hosts

## Vagrant Advantages

- reproducible and portable work environments
- easy to set up and learn
- usable for scale testing and development
- many *providers* available (Virtualbox, KVM, VMware...) to virtualize hosts
- choice between many *provisioners* (Shell, Ansible, Chef, Puppet...) to configure hosts

# Vagrant Advantages

- reproducible and portable work environments
- easy to set up and learn
- usable for scale testing and development
- many *providers* available (Virtualbox, KVM, VMware...) to virtualize hosts
- choice between many *provisioners* (Shell, Ansible, Chef, Puppet...) to configure hosts

# Example Vagrant Environment

- hardware with >8 cores and >32 GB RAM
- capable of hosting all OpenStack controller hosts full-scale
- Vagrant provider Virtualbox
- most-used provisioners: Shell, Ansible and Puppet

# Example Vagrant Environment

- hardware with >8 cores and >32 GB RAM
- capable of hosting all OpenStack controller hosts full-scale
- Vagrant provider Virtualbox
- most-used provisioners: Shell, Ansible and Puppet

# Example Vagrant Environment

- hardware with >8 cores and >32 GB RAM
- capable of hosting all OpenStack controller hosts full-scale
- Vagrant provider Virtualbox
- most-used provisioners: Shell, Ansible and Puppet

# Example Vagrant Environment

- hardware with >8 cores and >32 GB RAM
- capable of hosting all OpenStack controller hosts full-scale
- Vagrant provider Virtualbox
- most-used provisioners: Shell, Ansible and Puppet

# Vagrant at Work...

# Running Environment

# OpenStack Rally



Source: wpaphotomotorsport/Pixabay.com

# OpenStack Rally Advantages

- easy usage and setup
- many benchmark templates available which already cover many standard situations
- Rally plugins enable easy creation of more complex and use-case-specific benchmarks
- nice presentation of results

# OpenStack Rally Advantages

- easy usage and setup
- many benchmark templates available which already cover many standard situations
- Rally plugins enable easy creation of more complex and use-case-specific benchmarks
- nice presentation of results

# OpenStack Rally Advantages

- easy usage and setup
- many benchmark templates available which already cover many standard situations
- Rally plugins enable easy creation of more complex and use-case-specific benchmarks
- nice presentation of results

# OpenStack Rally Advantages

- easy usage and setup
- many benchmark templates available which already cover many standard situations
- Rally plugins enable easy creation of more complex and use-case-specific benchmarks
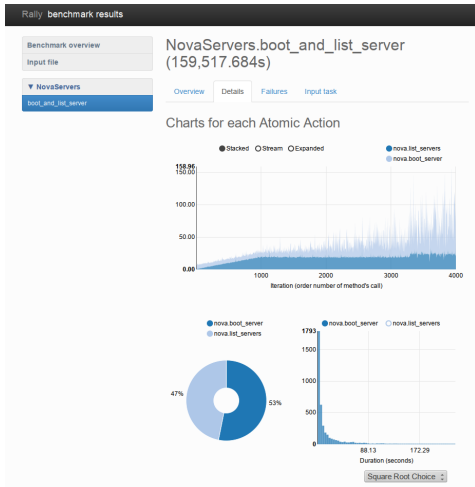- nice presentation of results

# OpenStack Rally in Action...

# OpenStack Rally

# Fake Drivers

- simulate instances or volumes
- transparent for the OpenStack controller hosts
- independent from hardware requirements

# Fake Drivers

- simulate instances or volumes
- transparent for the OpenStack controller hosts
- independent from hardware requirements

# Fake Drivers

- simulate instances or volumes
- transparent for the OpenStack controller hosts
- independent from hardware requirements

# Nova Fake Driver Configuration

- Fake Nova Compute Driver

## nova.conf

```
...
# Driver to use for controlling virtualization. Options
# include: libvirt.LibvirtDriver, xenapi.XenAPIDriver,
# fake.FakeDriver, baremetal.BareMetalDriver,
# vmwareapi.VMwareVCDriver, hyperv.HyperVDriver (string value)
compute_driver=fake.FakeDriver
...
```

# Conclusion

- determine clear design specifications (max instances, volumes, elasticity, users, tenants)
- use Rally to thoroughly test your setup within the specs
- perform a full-scale test without FakeDrivers prior to go-live
- use active/active HA setups for the core services to retain horizontal scalability

- determine clear design specifications (max instances, volumes, elasticity, users, tenants)
- use Rally to thoroughly test your setup within the specs
- perform a full-scale test without FakeDrivers prior to go-live
- use active/active HA setups for the core services to retain horizontal scalability

- determine clear design specifications (max instances, volumes, elasticity, users, tenants)
- use Rally to thoroughly test your setup within the specs
- perform a full-scale test without FakeDrivers prior to go-live
- use active/active HA setups for the core services to retain horizontal scalability

- determine clear design specifications (max instances, volumes, elasticity, users, tenants)
- use Rally to thoroughly test your setup within the specs
- perform a full-scale test without FakeDrivers prior to go-live
- use active/active HA setups for the core services to retain horizontal scalability

# Thank You!

For more information, refer to info@b1-systems.de
or +49 (0)8457 - 931096