

#### Postgres on OpenStack

• Dave Page | 18/9/2014



#### Introduction

- PostgreSQL:
  - Core team member
  - pgAdmin lead developer
  - Web/sysadmin teams
  - PGCAC/PGEU board member
- EDB:
  - Tools
  - Cloud products
  - Configuration management
  - Packaging/distribution
  - Infrastructure



#### What is OpenStack?

- Cloud environment, similar to AWS:
  - Private implementations
  - Public implementations Rackspace, HP Cloud
- Modular design, offering a choice of services including:
  - Image Service (Glance)
  - Block Storage (Cinder)
  - Compute (Nova)
  - Network (Neutron)
  - Object Storage (Swift)



#### **Glance Image Service**

- Store virtual machine images and snapshots
- Preconfigure an image or snapshot with the desired OS/software
- Share images with members of selected projects (AKA, tenants).



#### **Cinder Block Storage Service**

- Analogous to Amazon Elastic Block Storage (EBS)
- Provides block storage devices to virtual machine instances
- Can be provided from one or more service nodes via a controller
- Service nodes may use various storage systems, with different properties including:
  - LVM
  - GlusterFS
  - NFS
  - Various SANs



#### Nova Compute Service

- Analogous to Amazon Elastic Compute Service (EC2)
- Schedules and runs VM instances
- Multiple machine flavours with individual RAM, CPU and disk requirements
- Compute node selection is highly configurable:
  - Pluggable scheduler filters
  - Target images to specific hosts
  - Target flavours to specific hosts
  - Target instances to hosts with enough free resources



#### **Neutron Networking Service**

- Provides network services to VM instances
- Attach instances to one or more virtual networks
- Map floating/public IP addresses to a VM
- Configure firewall rules for ingress and egress



#### Swift Object Storage

- Analogous to Amazon Simple Storage Service (S3)
- Provides highly redundant object storage
- Suitable for storage and retrieval of whole files



#### Wait, what about Trove?

- Still in relatively early stages of development
- Currently seems focused on MySQL/Cassandra
- No support for clustering or other non-trivial setups
- API favours psql/pgAdmin operations, not deployment/ configuration



# Well that all sounds pretty cool, but what can we do with it?



## **Database Server Images**



#### Database server images

- Preconfigure a "standard" database image
- Enables consistency, which leads to ease of management
- Preinstall commonly required software:
  - Procedural languages
  - WAL-E
  - pgBouncer/pgPool
- Minimise launch/deployment time



#### Image or Snapshot?

- Technically there's not much difference
- OCD: I like images to be, well images
- Images may be considered "blessed" or approved
- Snapshots are a convenient way to take backups, perhaps of work in progress
- Images are "clean" configurations, whilst snapshots may contain artifacts from previous use, e.g. shell history or log files.



#### Creating an image

- Create a small root disk in QCOW2 format
- Install the base OS on the disk (I use KVM on RHEL)
- Install the required software
- Add the EPEL Yum repo, and install cloud-init
- Disable zeroconf for the network
- Shutdown the VM
- Run *virt-sysprep* to clean the OS installation



#### Uploading the image

\$	glance image-creat container-format	cename "PostgreSQL-9.4"disk-format qcow2 ∖ c bareis-public Trueprogress < pg-9.3.img
	Property	Value
	checksum   container_format   created_at   deleted_at   deleted_at   disk_format   id is_public   min_disk min_ram   name   owner   protected   size   status   updated_at	64d7c1cd2b6f60642c14663941cb7913 bare 2014-08-29T14:55:08 False None qcow2 acafc7c0-40aa-4026-9673-b879124bcec2 True 0 0 PostgreSQL-9.4 efa984b0a914450e9a47788ad330699d False 2013167616 active 2014-08-29T14:55:08
- T -		



## Database Server Machine Configurations



#### Dedicated database flavours

- "Flavors" (without the u) describe the instance configuration:
  - vCPUs
  - Memory
  - Root disk
  - Ephemeral disk
  - Swap
  - Other attributes
- Additional attributes allow us to target specific hypervisor properties, e.g.
  - Compute nodes with SSDs



## Targeting a host aggregate

- Create a host aggregate or group containing the desired compute nodes
- Add a property to the host aggregate to identify the desired attribute for matching
- Create the desired flavour(s)
- Add a property to the flavour to correspond with the property on the host aggregate
- Share the flavour with one or more tenants



#### Create host aggregate

#### \$ nova aggregate-create pg-servers +---+ | Id | Name | Availability Zone | Hosts | Metadata | +---+ | 14 | pg-servers | - | | | +---+

#### \$ nova aggregate-add-host pg-servers nova5.ox.uk Host nova5.ox.uk has been successfully added for aggregate 14 +---+ | Id | Name | Availability Zone | Hosts | Metadata | +---+ | 14 | pg-servers | - | 'nova5.ox.uk' | | +---+

#### \$ nova aggregate-set-metadata pg-servers pgsvr=true Metadata has been successfully updated for aggregate 14.

++	+		++
Id   Name	Availability Zone	Hosts	Metadata
14   pg-servers	-	'nova5.ox.uk' -+	'pgsvr=true'   ++



#### Create flavour



#### Create flavour

\$ nova flavor-key pgsvr.small set aggregate instance extra specs:pgsvr=true \$ nova flavor-show pgsvr.small Property | Value OS-FLV-DISABLED:disabled | False OS-FLV-EXT-DATA:ephemeral 0 disk 40 {"aggregate instance extra specs:pgsvr": "true"} extra specs 44221572-a7e0-4660-97fc-352467743fce id pgsvr.small name os-flavor-access:is public | True 4096 ram rxtx\_factor 1.0 4096 swap vcpus 4



#### A note on Scheduler Filters

- Scheduler filters determine which compute node a new instance will be created on
- Filters are chained, and can reject hosts based on:
  - Capabilities
  - Current utilisation
  - Physical location
  - Affinity
  - Other factors such as tenant isolation requirements
- Our example used the aggregate\_instance\_extra\_specs filter
- Custom filters can be defined in JSON (*json\_filter*), or written from scratch in Python



## Instance Storage -or-Where do my VMs live?



#### Local storage

- Known as *Instance Store* on AWS
- Images are copied to local storage on the compute node and started
- Pros:
  - Easy to setup
  - Performant
  - Minimal hardware requirements
- Cons:
  - Moving instances requires downtime
  - Potentially less reliable



#### Shared instance storage

- Same as Instance Store, but uses storage shared between compute nodes using NFS, iSCSI, etc.
- Pros
  - Allows easy migration of instances between compute nodes
  - Potentially more reliable storage (if a SAN or similar is used)
- Cons
  - "Local" I/O is over the network
  - More complex to setup and maintain
  - Can be far more expensive (SAN, 10GBe or fibre channel network etc.)



#### Volume storage

- Known as *Elastic Block Store* (EBS) on AWS
- Root volume is stored in Cinder
- Pros:
  - Available as a launch-time option, if Cinder is configured
  - Allows easy migration of instances between compute nodes
  - Potentially more reliable storage (if a SAN or similar is used)
- Cons:
  - "Local" I/O is over the network
  - More complex to setup and maintain
  - Can be far more expensive (SAN, 10GBe or fibre channel network etc.)



## Data Storage



#### Data storage

- In the instance:
  - Easy to setup and manage click and go
  - Resizing will pause the instance
  - Not very "elastic"
  - Slower snapshots data has to be moved to Glance
- In Cinder:
  - More complex to setup
  - More expensive for a high performance setup
  - Add storage on the fly by adding volumes and utilising LVM
  - Can choose storage types SSD vs. rust, low vs. high IOPs
  - Fast snapshots Cinder can use filesystem/SAN snapshots



# Firewalling



#### Instance firewalling

- Neutron provides Security Groups to define firewall
   rules
- No ingress access by default
- Each security group has one or more ingress or egress rules defined
- All instances have one or more security groups associated with them
- Allows easy configuration of a standard set of firewall rules for multiple servers



## Security group configuration

⊖   ⊖   ⊖									
🗲 🔿 C n 🍙 https://openstack-uk.enterprisedb.com/dashboard/project/access_and_security/security_group ☆ ı 🕑 👁 🥝 🚍									
dave.page 1									
Project Manage Security Group Rules: pg-svr									
Compute -	Sec	urity Gr	oup Rule	+ Add Rule					
Overview		Direction	Ether Type	IP Protocol	Port Range	Remote	Actions		
Instances		Ingress	IPv4	TCP	5432	0.0.0.0/0 (CIDR)	Delete Rule		
Volumes		Ingress	IPv4	ТСР	22 (SSH)	0.0.0.0/0 (CIDR)	Delete Rule		
Images		Egress	IPv4	Any	-	0.0.0.0/0 (CIDR)	Delete Rule		
Access & Security		Egress	IPv6	Any	-	::/0 (CIDR)	Delete Rule		
Network		Ingress	IPv4	ICMP	-	0.0.0.0/0 (CIDR)	Delete Rule		
Object Store	Displa	ying 5 items							
Admin 🕨									



# Networking



#### Private cluster network

- Create one or more virtual networks to contain a cluster of database servers, entire application or some other group of instances
- Analogous to Virtual Private Cloud (VPC) in Amazon
- Optionally attach virtual routers to networks
- Attach one or more networks to each instance
- Use *Floating IPs* (Elastic IPs) to access instances from outside the virtual network
- Each network can have a DHCP server and custom routes



## Network configuration

Image: Second state of the se										
← → C f la https://openstack-uk.enterprisedb.com/dashboard/project/networks/										
Infrastructure     ✓     Sign Out										
Project - Networks										
	Compute	•	Networks							
	Network	-						Admin		
	Network Topology			Name	Subnets Associated	Shared	Status	State	Actions	
	Networks			PG Cluster	PG Cluster 172.16.16.0/24	No	ACTIVE	UP	Edit Network More 👻	
	Routers			General VM Network	General VM Subnet 192.168.1.0/24	Yes	ACTIVE	UP		
	Object Store	•	Displaying 2 items							
A	Admin >									



#### Network topology





# WAL Archiving



## WAL archiving to Swift

- Object storage is ideal for storing WAL segments and backups
- Highly redundant storage offers reliability and scaling
- Easy to use with WAL-E (<u>https://github.com/wal-e/wal-e</u>)
  - WAL archival and retrieval
  - Base backup creation and restoration
  - WAL cleanup/pruning



## Installing WAL-E

```
# pip install wal-e
<stuff happens>
```

```
# vi /var/lib/pgsql/9.3/data/postgresql.conf
wal_level = archive
archive_mode = on
archive_command = '/var/lib/pgsql/wal-e.sh %p'
archive_timeout = 60 # For testing
```

#### # vi /var/lib/pgsql/wal-e-creds.sh

```
export SWIFT_AUTHURL=http://auth.enterprisedb.com:35357/v2.0
export SWIFT_USER=wal-e
```

```
export SWIFT_PASSWORD=ReallyReallySecret
```

```
export SWIFT_TENANT=admin
```

```
export WALE_SWIFT_PREFIX=swift://wal-e
```

```
# chown postgres:postgres /var/lib/pgsql/wal-e-creds.sh
# chmod 600 /var/lib/pgsql/wal-e-creds.sh
```



## Installing WAL-E

# vi /var/lib/pgsql/wal-e.sh #!/bin/sh source /var/lib/pgsql/wal-e-creds.sh /usr/bin/wal-e wal-push \$1 # chmod +x /var/lib/pgsql/wal-e-creds.sh # service postgresgl-9.3 restart # tail -f /var/lib/pgsql/9.3/data/pg log/postgresql.log MSG: starting WAL-E wal e.main INFO DETAIL: The subcommand is "wal-push". STRUCTURED: time=2014-08-29T16:55:56.236552-00 pid=18301 wal e.worker.upload INFO MSG: begin archiving a file STRUCTURED: time=2014-08-29T16:55:56.360962-00 pid=18301 action=push-wal state=begin urllib3.connectionpool INFO Starting new HTTP connection (1): auth.enterprisedb.com urllib3.connectionpool INFO Starting new HTTP connection (1): swift.enterprisedb.com wal e.worker.upload INFO MSG: completed archiving to a file at 9.29666KiB/s. STRUCTURED: time=2014-08-29T16:56:03.467501-00 pid=18301 action=push-wal key=swift://wal-e/wal 005/0000000100000000000000000000E.lzo prefix= rate=9.29666



WAL Archives

#### Containers

Containers + Create Container		Obj	Create Pseudo	Q Filter + Create Pseudo-folder			
wal- e	Object Count: 17 Size: 13.4 MB Access: Private	View Details More *		extended_version.txt	108 bytes	Download More *	
				part_0000000.tar.lzo	6.9 MB	Download More *	
				base_0000000100000000000000000000000000000	202 bytes	Download More *	
				00000010000000000000000000000000000000	65.9 KB	Download More *	
				00000010000000000008.lzo	65.9 KB	Download More *	
				0000001000000000000009.lzo	65.9 KB	Download More *	
				0000001000000000000A.lzo	65.9 KB	Download More *	
				0000001000000000000B.lzo	2.7 MB	Download More *	
				000000010000000000000000C.lzo	65.9 KB	Download More *	



# Summary



#### Summary

- Deploy standardised, maintainable database servers from images in Glance
- Configure database-optimised compute nodes
- Configure database-optimised machine flavours in Nova to target desired nodes
- Select instance storage based on performance and flexibility requirements vs. budget.
- Select data storage based on complexity vs. flexibility



#### Summary contd.

- Use Neutron's security groups to provide standardised security configuration outside of the instance
- Utilise virtual networks to isolate logical clusters or groups of machines
- Archive base backups and WAL to Swift with WAL-E for ease of backup management, redundancy and scalability



## Questions?





