



TECH TALK:

Cost Estimation and Control for Researchers

Friday, December 4th, 2020
10:30 AM US EDT

Rajib Ghosh
Global Senior Solutions Architect
Oracle for Research

TECH TALK HOUSEKEEPING

- Today's webinar is being recorded. We will share the link to the recording with you via email after the event. The recording will also be made available to the Oracle for Research community.
- We invite your comments and questions, both about the tech topic being discussed and about the series more generally. Questions may be submitted using the Q&A box on your screen or you may ask questions directly using your microphone. When not asking a question, please mute your microphone.
- Questions may be asked during the presentation and we will also have a Q & A time at the end of the presentation when you can ask questions directly and engage in discussion.
- At Oracle for Research, we believe that research and innovation happen best when a diverse and thoughtful community is free to engage in respectful, compassionate, and open dialog. To that end, when asking a question or providing feedback, we ask that all participants be respectful, collaborative, and constructive.

Agenda

Recap and Asks from researchers

1. Cost analysis and estimation
2. How to prevent run-away tenancy costs?
3. Any recommendations for cost control?

Cost Management

1. Cloud advisor
2. Guidelines for researchers

Cost estimation

1. Estimation tools and process
2. Guidelines for researchers

Cost analysis

1. Cost analysis and reporting
2. Guideline for researchers

Q & A

Resources
Oracle for Research Github collaboration

Key factors of importance

For Researchers

- ❖ Minimizes idle time for cloud resources
- ❖ More computational cycles per credit \$
- ❖ Preserve cloud credit for large computations
- ❖ Analyze potential cost over-run areas
- ❖ Prevent credit drainage
- ❖ Reduced resource contention
- ❖ Better control on projects / workloads

For Oracle

- ❖ Better utilization of cloud resources
- ❖ Helps sustain more research projects
- ❖ Gain knowledge on research computing usage
- ❖ Provide improved cost management tools
- ❖ Better cost analytics and reports for research
- ❖ Better monitoring of cloud compute shapes
- ❖ Quicker provisioning of service limit requests

Cloud advisor

Overview

1. Automated guidance that finds potential tenancy inefficiencies
2. Provides recommendations and cost savings in \$\$
3. Built into OCI platform
4. OCI CLI enabled

Key benefits

1. Downsize underutilized compute instances
2. Resize underutilized Autonomous databases and instances
3. Manages orphaned block and boot volumes
4. Customized recommendations with monitoring enablement
5. Migrate data with low cost storage based on lifecycle policy rules

Recommendations and Cost calculations

1. Based on previous month's data or month-to-date data (newer objects)
2. Compute recommendations – Based on trailing 7 days of data
3. Compute costs = (billed usage * unit price) / 2
4. Block / boot volume recommendations – based on billed usage, performance units and storage

Researcher Guidelines

1. Check if monitoring for CPU / GPU instances. (Default for all Oracle images or Custom images built on it)
2. Downsize VM shapes / databases based on recommendations
3. Build instances from custom images instead of boot volumes. Always delete boot volumes
4. Enable Object store lifecycle management (OLM) for long term data archival and deletion
5. Customize recommendation profiles based on your usage and utilization requirements
6. Start with average (default) methodology and configure P95 if required
7. Use one compartment / project and shared resources in separate compartments

Cost estimation tools for research

Tools

- ❖ Infrastructure cost pricing - <https://www.oracle.com/cloud/price-list.html#compute-gpu>
- ❖ Oracle cloud cost estimator - <https://www.oracle.com/cloud/cost-estimator.html>
- ❖ Oracle cloud workload estimator - <https://www.oracle.com/webfolder/workload-estimator/index.html>

Compute - GPU Instances

Instances are available as both virtual machines and bare metal, providing flexibility and performance at the fraction of the cost of other public cloud

Shape	GPUs	Architecture	GPU Interconnect	GPU Memory	CPU Cores	CPU Memory
VM.GPU2.1	1x NVIDIA P100	Pascal	N/A	16 GB	12	78 GB
BM.GPU2.2	2x NVIDIA P100	Pascal	N/A	32 GB	28	192 GB
VM.GPU3.1	1x NVIDIA V100 Tensor Core	Volta	N/A	16 GB	6	90 GB
VM.GPU3.2	2x NVIDIA V100 Tensor Core	Volta	NVIDIA NVLINK	32 GB	12	180 GB
VM.GPU3.4	4x NVIDIA V100 Tensor Core	Volta	NVIDIA NVLINK	64 GB	24	360 GB
BM.GPU3.8	8x NVIDIA V100 Tensor Core	Volta	NVIDIA NVLINK	128 GB	52	768 GB
VM.GPU4.1*	1x NVIDIA A100 Tensor Core	Ampere	N/A	40 GB	7	224 GB
VM.GPU4.2*	2x NVIDIA A100 Tensor Core	Ampere	NVIDIA NVLINK	80 GB	15	480 GB
VM.GPU4.4*	4x NVIDIA A100 Tensor Core	Ampere	NVIDIA NVLINK	160 GB	30	960 GB
BM.GPU4.8	8x NVIDIA A100 Tensor Core	Ampere	NVIDIA NVLINK	320 GB	64	2048 GB

*Available soon

Utilization

▶ Number of Instances / 1 Instance(s)

▲ Average Days Usage per Month / 16 day(s)

Indicate average days usage per month for this service

▲ Average Hours Usage per Day / 13 hour(s)

Indicate average hours usage per day for this service

Configuration

▶ Compute - Virtual Machine Standard - X7 (B88514) / 1 OCPU Per Hour	\$13
▶ Compute - BM Standard - B1 (B91119) / 1 OCPU Per Hour	\$13
▶ Compute - HPC - X7 (B90398) / 1 OCPU Per Hour	\$16
▶ Compute - Microsoft SQL Enterprise - OCPU Per Hour (B91372)	\$0
▶ Compute - Microsoft SQL Standard - OCPU Per Hour (B91373)	\$0
▶ Compute - Standard - E2 (B90425) / 1 OCPU Per Hour	\$6
▶ Compute - VM Standard - B1 (B91120) / 1 OCPU Per Hour	\$13
▶ Compute - Virtual Machine Dense I/O - X7 (B88516) / 1 OCPU Per Hour	\$27
▶ Compute - Virtual Machine GPU Standard - X7 (B88518) / 1 GPU Per Hour	\$265
▶ Oracle Cloud Infrastructure - Compute - GPU - E3 - GPU Per Hour (B92740) / 1 GPU Per Hour	\$634
▶ Virtual Machine Standard - X5 (B88317) / 1 OCPU Per Hour	\$13
▶ Windows OS (B88318)	\$0

Cost estimation

GPU / CPU hours and Cost estimation

STEP-1: (Researcher On-campus estimation)

- ❖ Record test computation hours on campus dedicated machine (laptop/physical server)
- ❖ Compute total estimated GPU/CPU hours
- ❖ Record hardware details (CPU/GPU/RAM, spec, storage)
- ❖ Cloud bursting ratio (if applicable)
- ❖ Provide details to Oracle for Research

STEP-2: (Oracle cloud self-service estimation)

- ❖ Compute GPU / CPU hours / hour
- ❖ Estimate GPU / CPU cores / hour
- ❖ Choose nearest Bare metal shape
- ❖ Compute required cluster nodes / shape
- ❖ Compute spend \$/node and \$/month
- ❖ Adjust \$ spend based on bursting ratio / execution frequency
- ❖ Provide details to Oracle for Research

STEP-3: Compute and request for service limit increase

Cost estimation – Example

Science Driver	CPU	GPU Khrs/yr	kHrs/Month	kHrs/hr	GPU cores	Storage			
Computer vision	Minimal	200	16.67	23.15	24	100TB			
Products	Shape	Specification	Part #	Rate	Metric	\$/Node	\$max/month	Quantity/ Service limits	Assumptions / Comments
Compute (V100 GPU)	OCI VM.GPU3.1	1xGPU + 16GB-GPU-RAM + 6xoCPU + 90GB-CPU-RAM	B89734	2.9500	GPU / hr	566.42	13594	24	
	OCI VM.GPU3.2	2xGPU + 32GB-GPU-RAM + 12xoCPU + 180GB-CPU-RAM	B89734	2.9500	GPU / hr	1,132.83	13594	12	
	OCI VM.GPU3.4	4xGPU + 64GB-GPU-RAM +24xoCPU + 360GB-CPU-RAM	B89734	2.9500	GPU / hr	2,265.67	13594	6	
	OCI BM.GPU3.8 (V100)	8xGPU + 128GB-GPU-RAM +52xoCPU + 768GB-CPU-RAM	B89734	2.9500	GPU/hr	4,531.33	13594	3	
	OCI.BM.GPU4.8 (A100)	8xGPU + 320GB-GPU-RAM +64xoCPU + 2048GB-CPU-RAM + 8x200 Gbps RDMA		3.0500	GPU/hr	4,685.00	14055	3	
Block volume	Storage		B91961	0.0255	GB/Month		2550	100TB	
	Performance units	Balanced	B91962	0.0170	GB/month		1700	100TB	Default performance unit
		High Performance	B91962	0.0340	GB/Month		3400	100TB	Based on full BV capacity
						Total (Max)	19544		Compute + BV storage + BV (high Perf)

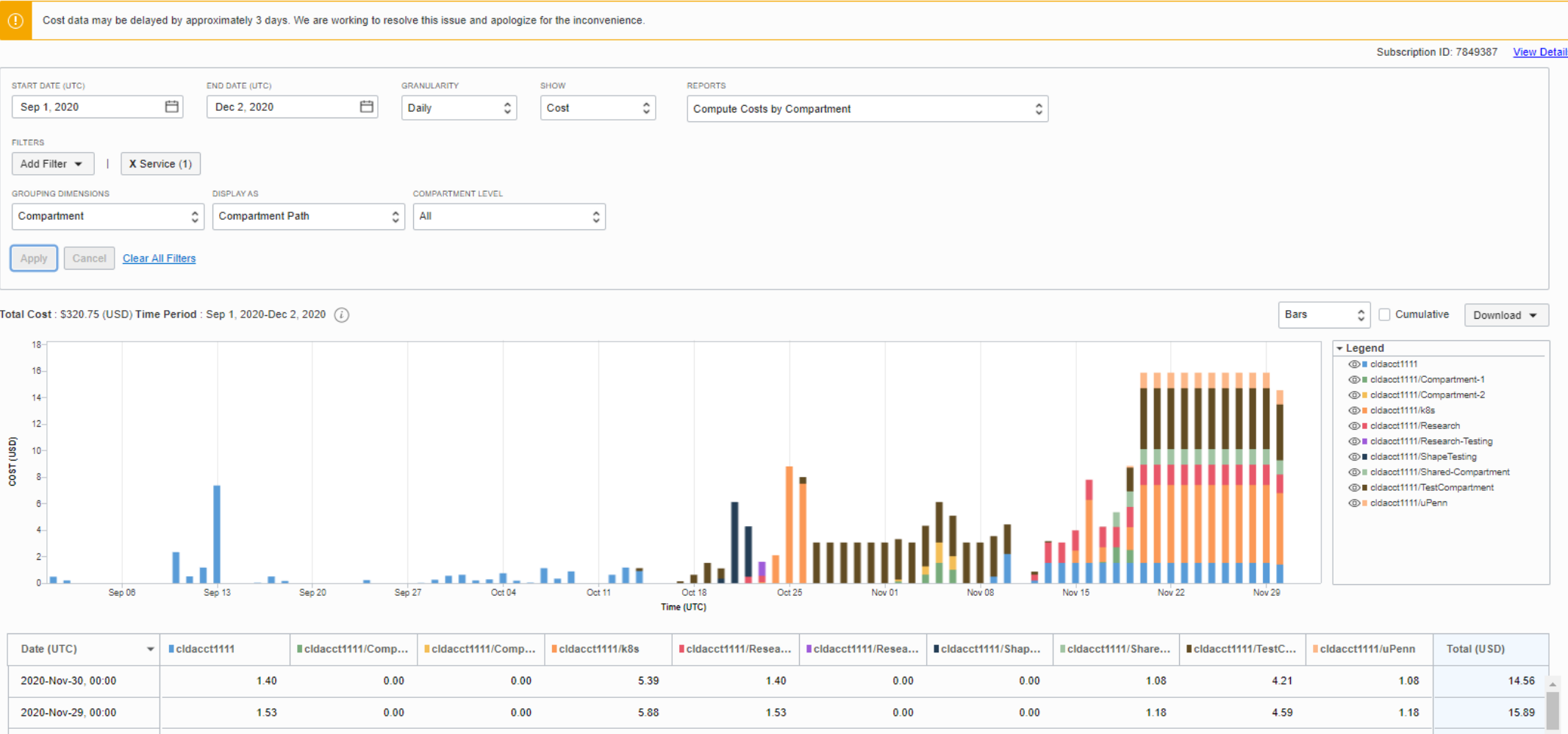
1. Actual GPU hours for BM machines can be less as they are dedicated machine with no overhead.
2. Actual GPU hours can be much less for BMGPU4.8 (320GB GPU / better Tensor flow architecture)

Typical cost analysis for research

Reports, filter and dimensions

Report: Compute costs by compartment
Filter: Service (Compute + block storage)
Dimension: compartment All
Optional filters: Granularity (monthly/Daily) & show (cost/usage)

Cost Analysis



Cost analysis for research - Advanced

Using Tags

Report: Compute monthly costs by compartment
Filter: Service (Compute + block storage) + Oracle Tag (Created by)
Dimension: compartment (All)
Optional filters: Show (cost/usage)

Cost Analysis

Cost data may be delayed by approximately 3 days. We are working to resolve this issue and apologize for the inconvenience.

Subscription ID: 7849387 [View Details](#)

START DATE (UTC)
Aug 1, 2020

END DATE (UTC)
Dec 2, 2020

GRANULARITY
Monthly

SHOW
Cost

REPORTS
--

FILTERS

Add Filter

X Service (1)

X Tag: Oracle-Tags.CreatedBy : oracleidentitycloudservice/ra_rghosh@yahoo.com

GROUPING DIMENSIONS
Compartment

DISPLAY AS
Compartment Path

COMPARTMENT LEVEL
All

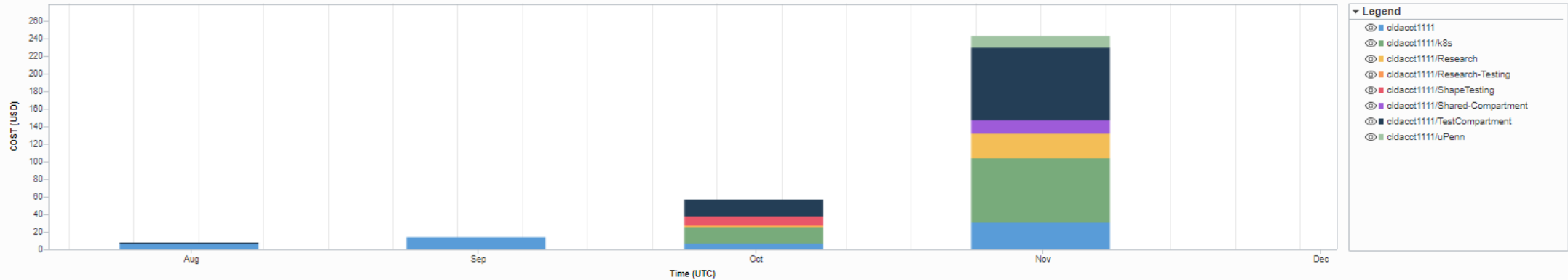
Apply

Cancel

[Clear All Filters](#)

Total Cost : \$320.59 (USD) Time Period : Aug 1, 2020-Dec 2, 2020

Bars Cumulative Download



Date (UTC)	cldacct1111	cldacct1111/k8s	cldacct1111/Research	cldacct1111/Research-Testing	cldacct1111/ShapeTesting	cldacct1111/Shared-Compartment	cldacct1111/TestCompartment	cldacct1111/uPenn	Total (USD)
2020-Nov	30.44	73.19	28.00	0.00	0.00	15.23	82.46	13.06	242.39
2020-Oct	6.78	18.48	0.96	1.06	10.26	0.00	19.11	0.00	56.65



Cost analysis for research

Guidelines

Simple use-case - (Quick cost analysis by a PI)

- ❖ Allocate a compartment for a project
- ❖ Track resource usage and cost by compartment

Complex use-case – (Multiple users and projects)

- ❖ Use Oracle provided tagging or custom tagging
- ❖ Tag Oracle cloud instance at creation
- ❖ Track resource usage/cost with tags
- ❖ Use OCI CLI usage API for automation
- ❖ Download cost and usage reports for details

Setting Budgets

- ❖ Set budget alert at root and heavy usage compartments.
- ❖ Configure PI + researcher email for budgeting notifications

Reporting

- ❖ Download and customize automatic detail csv reports for further analysis
- ❖ Keep only relevant columns for analysis
- ❖ Feed into downstream analytics systems
- ❖ May use OCI usage API for automation

Oracle cloud Cost management links

Frequently accessed links for research

[Oracle cloud price list and cost estimator](#)

[Resource billing for stopped instances](#)

[Universal credit for PaaS and IaaS services](#)

Github and documentation

[Python SDK to set up Usage and Cost Reports](#)

[OCI Reporting Tool](#)

[OCI Cost reports to Autonomous database tool](#)

OCI-CLI blogs

[OCI Cost reports with APEX app](#)

[OCI rate card utility](#)

[OCI General cloud forum](#)

[Using python and metering API to produce custom cost reports](#)

[Working with OCI cost reports](#)

[Tracking costs with OCI Tagging defaults](#)

[State of the art cost governance on OCI](#)

[Auto tune detached block volume to save cost](#)

[Building high value, low cost data lakes](#)

[Using quotas for effective cost management in cloud](#)

[What everyone should know about cloud consumption](#)

[Best practices for using Tags to manage costs, operations and governance](#)



TECH TALK:

Cost Estimation and Control for Researchers

Questions, Answers & Discussion



TECH TALK:

Cost Estimation and Control for Researchers

Questions? Comments? Feedback?

Contact us!

Website: oracle.com/oracle-for-research/

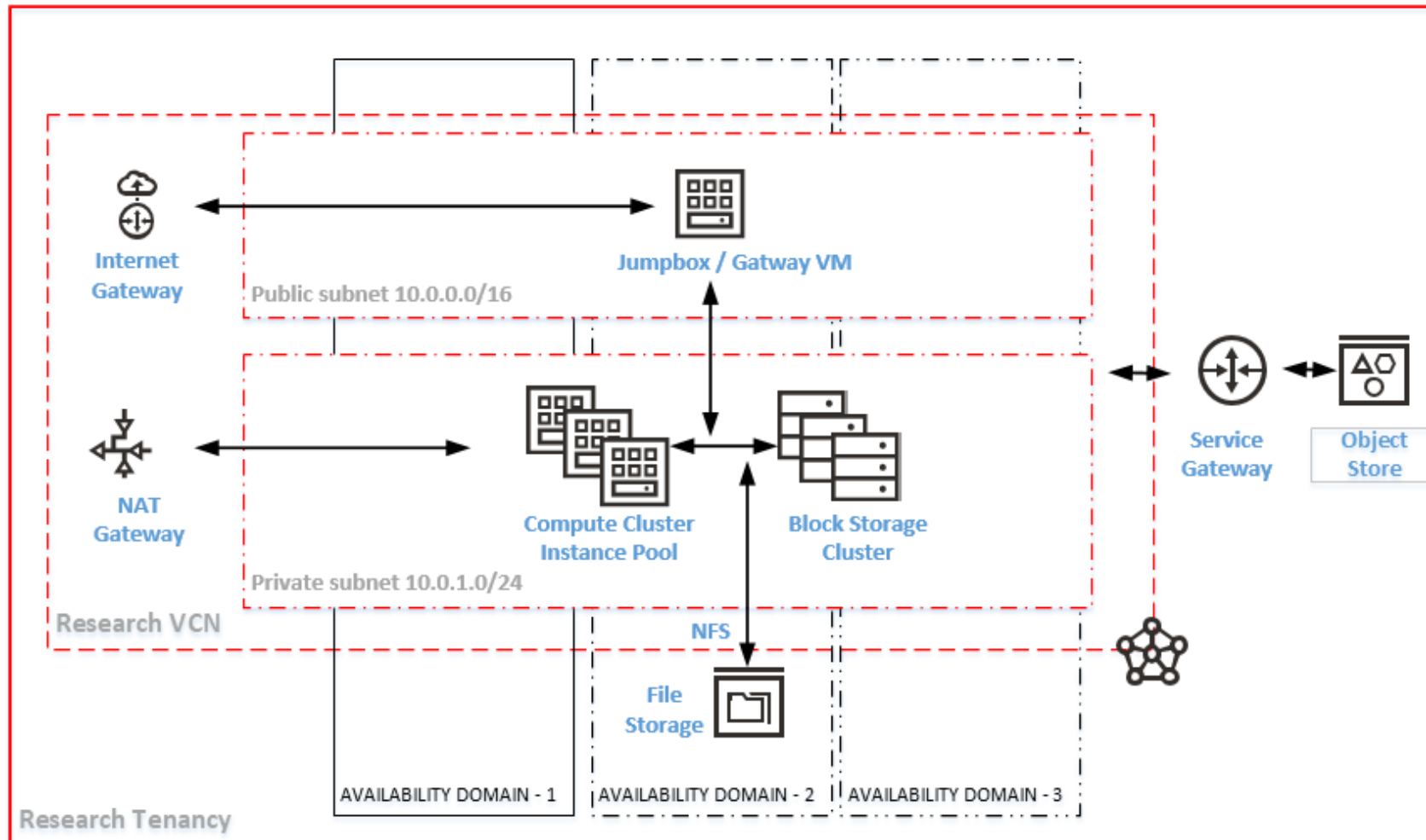
Github: github.com/OracleforResearch

Twitter: @OracleResearch

Email: OracleForResearchTech_ww@oracle.com

Next Tech Talk: Jan 2021 – date to be announced

OCI Standard cluster architecture for Researchers



Component	Recommendation
Jumpbox VM	Use Free-Tier VM
CPU/GPU Cluster	Manual Build + Block storage
CPU/GPU Cluster	Automated Build + Block storage
HPC Cluster	Instance Pool + RDMA
NAT Gateway	Download software to pvt subnet
File storage	For cross-AD exports

Instance Type	Shape series	Shape	Purpose
Virtual	Always Free	VM.StandardE2.1Micro	Automation control, gateway, configurations
	Standard	VM.Standard1.1~1.16	Low workload testing / Image builds / installs
	AMD (Gen 2)	VM.StandardE2.1~2.8	Prototype workload testing
	DenseIO	VM.DenseIO2.x (NVMe)	Heavy IO workload testing
	GPU (P100)	VM.GPU2.1	AI / ML or other GPU prototype testing
	GPU (V100)	VM.GPU3.1~3.4	Tensor core AI / DL workloads
	Intel Skylake (Fixed)	VM.Standard2.1~2.24	Workloads to save on credits
	AMD Rome (Flex)	VM.StandardE3.Flex	Benchmarking / price-performance
Bare metal	HPC	BM.HPC2.36 (NVMe)	CPU+high throughput for HPC workloads
	AMD (Gen 3)	BM.StandardE3.128	High CPU/throughput workloads
	Standard	BM.Standard1.36/B1.44	Low CPU/RAM utilization at lowest BM cost
	AMD (Gen 2)	BM.StandardE2.52	Best price-performance for BM workloads
	AMD (Gen 3)	BM.StandardE2.64	Best Gen3 price-performance for BM workloads
	DenseIO	BM.DenseIO2.52 (NVMe)	Best price performance for IO intensive workloads
	GPU (P100)	BM.GPU2.2	Benchmarking pascal based GPU workloads
	GPU (V100)	BM.GPU3.8	Best price performant for large GPU workloads
	GPU (A100)	BM.GPU4.8	Fastest GPU – large DL applications (pre-GA)