

# Вычислительные мощности Балансировщики нагрузки в Cloud Provider

Елисей Ильин



# План занятия

1. [Вычислительные мощности](#)
2. [Хранение](#)
3. [Балансировщики нагрузки](#)
4. [Итоги](#)
5. [Домашнее задание](#)



# **Вычислительные мощности (виртуальные машины)**

---

# AWS EC2 — Elastic Compute Cloud (ECC)

AWS EC2 типы тарификации:

- **On Demand** — объём вычислительных ресурсов оплачивается на почасовой или посекундной основе в зависимости от используемых инстансов
- **Saving Plans** — объём, предусматривающий снижение оплаты при условии, что клиент обязуется использовать этот постоянный объём в течение одного года или трёх лет
- **Spot** — объём свободных вычислительных ресурсов Amazon EC2 со скидкой до 90% по сравнению с ценой по требованию
- **Dedicated** — выделенный объём физических серверов, предоставляемый только одному клиенту



# Типы EC2-инстансов

Типы инстансов включают различные комбинации таких компонентов, как:

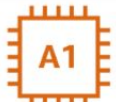





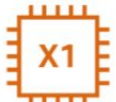



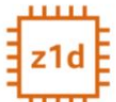
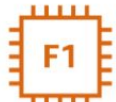

- ЦПУ
- память
- хранилище
- сетевые возможности

Это позволяет выбрать соответствующий набор ресурсов для приложений

Каждый тип инстанса включает **один** или **несколько размеров** инстансов, что позволяет масштабировать ресурсы в соответствии с требованиями целевой рабочей нагрузки

# Типы EC2-инстансов

- General Purpose
- Compute Optimised
- Memory Optimised
- Accelerated Computing
- Storage Optimised

General Purpose	Compute Optimised	Memory Optimised	Accelerated Computing	Storage Optimised
 ARM based core and custom silicon	 Compute - CPU intensive apps and DBs	 RAM - Memory intensive apps and DB's	 Processing optimised- Machine Learning	 High Disk Throughput - Big data clusters
 Tiny - Web servers and small DBs		 Xtreme RAM - For SAP/Spark	 Graphics Intensive - Video and streaming	 IOPS - NoSQL DBs
 Main - App servers and general purpose		 High Compute and High Memory - Gaming	 Field Programmable - Hardware acceleration	 Dense Storage - Data Warehousing

Источник

---

## Сетевые адаптеры для EC2

- **ENI** — (Elastic network interfaces) виртуальная сетевая карта инстанса
- **Enhanced network** — высокоскоростная сетевая карта 10 Гбит/с, 100 Гбит/с
- **Elastic Fabric Adapter** используется для HPC и ML, где высокие требования к задержкам

---

# YC Compute Cloud

- **Виртуальная машина** — объём вычислительных ресурсов оплачивается посекундно. Цены указаны за один час использования
- **Прерываемая виртуальная машина** — виртуальные машины, которые могут быть принудительно остановлены в любой момент

Это может произойти в двух случаях:

- с момента запуска виртуальной машины прошло 24 часа
- нехватка ресурсов для запуска обычной виртуальной машины в той же зоне доступности. Вероятность такого события низкая, но может меняться изо дня в день



# УС – платформы Compute Cloud

Платформа	Процессор	Макс. кол-во ядер (vCPU) на ВМ	Макс. кол-во GPU (vGPU*)
Intel Broadwell	<a href="#">Intel® Xeon® Processor E5-2660 v4</a>	32	-
Intel Cascade Lake	<a href="#">Intel Xeon Gold 6230</a>	80	-
Intel Ice Lake	<a href="#">Intel Xeon Gold 6338</a>	96	-
Intel Broadwell with NVIDIA® Tesla® V100	<a href="#">Intel Xeon Processor E5-2660 v4</a> <a href="#">NVIDIA® Tesla® V100</a>	32	4 (1:8)
Intel Cascade Lake with NVIDIA® Tesla® V100	<a href="#">Intel Xeon Gold 6230</a> <a href="#">NVIDIA® Tesla® V100</a>	64	8 (1:8)
AMD EPYC™ with NVIDIA® Ampere® A100	<a href="#">AMD EPYC™ 7702</a> <a href="#">NVIDIA® Ampere® A100</a>	224	8 (1:28)
Intel Broadwell with NVIDIA® vGPU Tesla® V100 8G	<a href="#">Intel Xeon Processor E5-2660 v4</a> <a href="#">NVIDIA® Tesla® V100</a>	32	8 (1:4)*

## УС — пример цен

	Вычислительные ресурсы ВМ	Вычислительные ресурсы прерываемых ВМ
Intel Broadwell, 5% vCPU	0,3100 ₽	0,1900 ₽
Intel Broadwell, 20% vCPU	0,8800 ₽	0,2700 ₽
Intel Broadwell, 100% vCPU	1,1200 ₽	0,3400 ₽
Intel Cascade Lake, 5% vCPU	0,1600 ₽	0,1000 ₽
Intel Cascade Lake, 20% vCPU	0,4900 ₽	0,1600 ₽
Intel Cascade Lake, 50% vCPU	0,7200 ₽	0,2200 ₽
Intel Cascade Lake, 100% vCPU	1,1900 ₽	0,3200 ₽

## УС — уровни производительности vCPU

Этот уровень определяет долю вычислительного времени физических ядер, которую гарантирует vCPU

- При уровне производительности 5% VM будет иметь доступ к физическим ядрам как минимум 5% времени — 50 миллисекунд в течение каждой секунды. Тактовая частота процессора в это время не ограничивается и соответствует выбранной платформе, например, 2 ГГц для платформы Intel Ice Lake (standard-v3)

VM с уровнем производительности меньше 100% предназначены для запуска приложений, не требующих высокой производительности и не чувствительных к задержкам. Такие машины обойдутся дешевле



## УС — уровни производительности vCPU

- Виртуальные машины с уровнем производительности 100% имеют непрерывный доступ (100% времени) к вычислительной мощности физических ядер

Такие ВМ предназначены для запуска приложений, требующих высокой производительности на протяжении всего времени работы



# Хранение

---

## AWS — типы Storage

- **EBS** — (Elastic Block Storage) — сервис блочного хранилища, созданный для использования с EC2
- **EFS** — (Elastic File Storage) сервис файлового хранилища, который позволяет **совместно** использовать файловые данные без необходимости его обслуживания
- **S3** — (Simple Storage Service) сервис хранения объектов (имеют ключ, значение, версию, метаданные, ACL). С помощью API обеспечивает доступ к данным через интернет

# Типы EBS

	SSD		HDD	
Наименование	General Purpose	Provisioned IOPS	Throughput Optimized	Cold
API Name	(gp2, gp3)	(io1, io2)	(st1)	(sc1)
Примеры использования	Часто используемые	База данных	Big Data, DWH	Файловые серверы
Емкость	1 GB – 16 TB	4 GB – 16 TB	125 GB – 16 TB	125 GB – 16 TB
Max IOPS*/Volume	16 000	64 000	500	250
Цена в месяц	\$0,08–0,1 /GB	\$0,125 /GB	\$0,045 /GB	\$0,015/GB

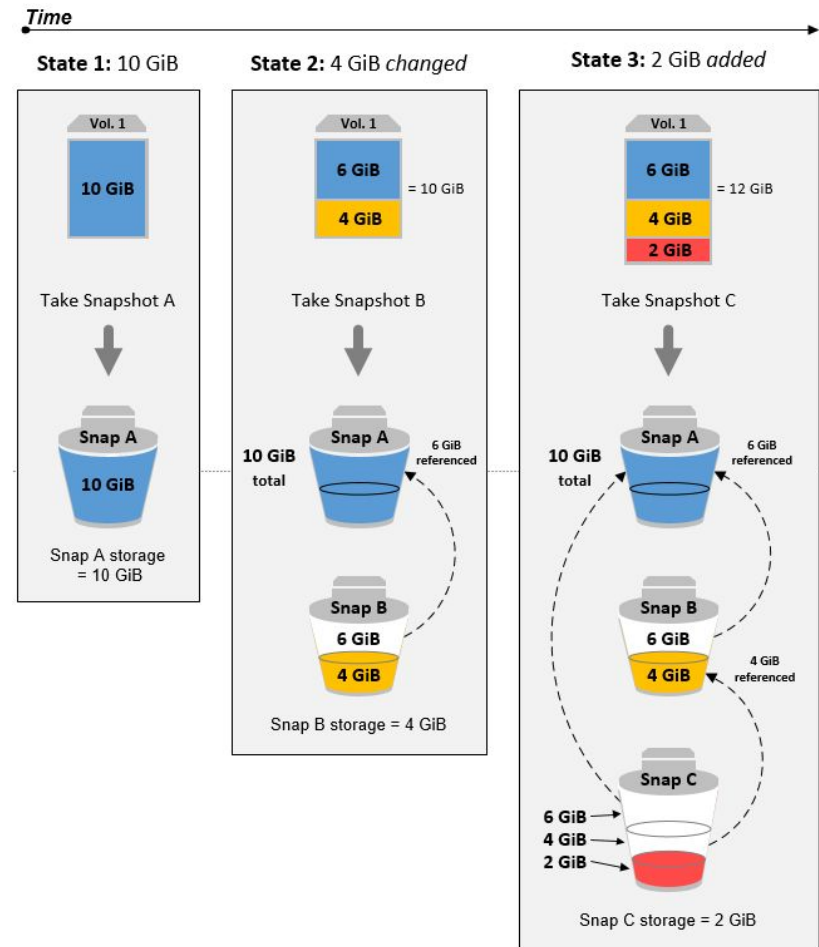
\* – io1/io2/gp2/gp3 based on 16K I/O size, st1/sc1 based on 1 MB I/O size

[Источник](#)

# EBS Snapshots

Моментальный снимок EBS, содержащий копии файлов и каталогов файловой системы на определённый момент времени:

- Snapshot хранить в S3
- Snapshot инкрементальный
- Поддержка шифрования
- Для создания рекомендуется остановить EC2 (консистентность)
- Используется для переноса EC2 из одного AZ или региона в другой с помощью AMI





---

# EFS

EFS предоставляет четыре класса хранилищ:

- **EFS Standard** — данные хранятся с избыточностью в разных AZ
- **EFS Standard-IA (Infrequent Access)** — нечастый доступ EFS Standard
- **EFS One Zone** — данные хранятся с избыточностью в одной AZ
- **EFS One Zone-IA** — нечастый доступ EFS One Zone

Особенности EFS:

- EFS поддерживает протокол NFSv4
- Не требуется расширять пространство, оплата по факту использования
- Поддерживает множество одновременных подключений
- EFS не поддерживает Windows instances



## S3

Данные хранятся, как **объекты** в ресурсах, которые называют корзинами (bucket), при этом размер одного объекта может составлять до 5 ТБ

Доступ к объектам можно получить через точки доступа S3 или непосредственно через имя узла контейнера

---

# S3

Хранилище S3 позволяет:

- добавлять **теги** метаданных в объекты
- перемещать и сохранять данные **в классах хранилища S3**
- настраивать элементы управления **доступом к данным**
- применять аналитику больших данных
- отслеживать данные на уровне объекта и корзины
- просматривать статистику использования хранилищ и тенденции активности в своей организации

## S3-типы доступа

- Virtual-hosted-style access

`https://bucket-name.s3.Region.amazonaws.com/file_name`

например:

`https://netology-devops15.s3.us-west-2.amazonaws.com/cat.png`

- Path-style access

`https://s3.Region.amazonaws.com/bucket-name/file_name`

например:

`https://s3.us-west-2.amazonaws.com/netology-devops15/cat.png`

---

## Классы хранилищ S3

S3 предоставляет несколько классов хранилищ:

- **S3 Standard** — данные хранятся с избыточностью в разных AZ
- **S3 Standard-IA** — нечастый доступ S3 Standard
- **S3 One Zone-IA** — нечастый доступ S3 Standard в одной AZ
- **S3 Intelligent-Tiering** — автоматическая оптимизация затрат путём переноса на разные классы хранилищ S3
- **S3 Glacier** — экономичный класс для архивации данных
- **S3 Glacier Deep Archive** — класс долгосрочного хранения и цифровой архивации данных, доступ запрашивается один–два раза в год

# Классы хранилищ S3

	S3 Standard	S3 Intelligent-Tiering	S3 Standard – IA	S3 One Zone-IA	S3 Glacier	S3 Glacier Deep Archive
Спроектировано для доступности	99,99%	99,9%	99,9%	99,5%	99,99%	99,99%
Доступность согласно SLA	99,9%	99%	99%	99%	99,9%	99,9%
Зоны доступности	≥3	≥ 3	≥ 3	1	≥ 3	≥3
Минимальный оплачиваемый объём объекта	Н/д	Н/д	128 КБ	128 КБ	40 КБ	40 КБ
Минимальный оплачиваемый срок хранения	Н/д	30 дней	30 дней	30 дней	90 дней	180 дней
Плата за извлечение данных	Н/д	Н/д	За гигабайт извлечённых данных	За гигабайт извлечённых данных	За гигабайт извлечённых данных	За гигабайт извлечённых данных
Задержка первого байта	Мсек	Мсек	Мсек	Мсек	Мин или часы на выбор	Выбрать часы

[Источник](#)

---

## УС — типы Storage

- **Диск** — сервис блочного хранилища, созданный для использования с Compute Cloud
- **Файловое хранилище** — сервис файлового хранилища, который позволяет **совместно** использовать файловые данные без необходимости его обслуживания
- **Object Storage** — сервис хранения объектов, с помощью API обеспечивает доступ к данным через интернет — т. е. из любой точки. Совместим с API S3, поэтому можно использовать инструменты, созданные для работы с объектными хранилищами

## Типы дисков УС

- **Сетевой SSD-диск** — быстрый сетевой диск, сетевое блочное хранилище на SSD-накопителе
- **Сетевой HDD-диск** — стандартный сетевой диск, сетевое блочное хранилище на HDD-накопителе
- **Нереплицируемый SSD-диск** — сетевой диск с повышенной производительностью, реализованной за счёт некоторых ограничений:
  - размер нереплицируемого диска должен быть кратен 93 ГБ
  - во всех расчётах 1 ГБ =  $2^{30}$  байт
  - нереплицируемые диски не могут быть загрузочными
  - хранимая информация может быть временно недоступна или утеряна в случае сбоя, т. к. в нереплицируемых дисках отсутствует избыточность
  - из нереплицируемого диска нельзя создать **снимки** и **образы**



---

## Репликация дисков и снимки

- Каждый диск доступен и реплицируется внутри определённой **зоны доступности**
- Можно создавать резервные копии дисков в виде **снимков**. Снимки реплицируются во всех зонах доступности. Можно использовать, чтобы переносить диски между зонами доступности
- Несколько **нереплицируемых** дисков могут быть собраны в группу размещения — группу нереплицируемых дисков, в которой диски располагаются в разных стойках в пределах одной зоны доступности
- Группы размещения используются для организации избыточности хранения данных на уровне приложения

---

## YC Object Storage

- Имена бакетов уникальны для всего Object Storage
- Длина имени должна быть от 3 до 63 символов
- Имя может содержать строчные буквы латинского алфавита, цифры, дефисы и точки

---

# Object Storage — типы доступа

- Virtual-hosted-style access

`https://bucket.storage.yandexcloud.net`

например:

`https://netology-test.storage.yandexcloud.net/cat.jpg`

- Path-style access

`https://storage.yandexcloud.net/bucket/file_name`

например:

`https://storage.yandexcloud.net/netology-test/cat.jpg`

# Классы хранилищ Object Storage

Object Storage предоставляет три класса хранилищ:

- **стандартное** — частые запросы
- **холодное** — редкие запросы
- **ледяное** — очень редкие запросы

Чем холоднее хранилище, тем дешевле хранить в нем данные, но тем дороже их читать и записывать.



# S3-совместимость Object Storage

**Yandex Object Storage** API совместим с **AWS S3** API

S3-совместимость позволяет использовать в Yandex Cloud популярные инструменты для работы с S3-протоколом:

- консольные клиенты S3cmd и AWS CLI
- файловые браузеры Cyberduck и WinSCP
- библиотеки (SDK) для Python и Java



# Балансировщики нагрузки



# AWS ELB — балансировщик нагрузки

**Elastic Load Balancer** автоматически распределяет входящий трафик приложений по нескольким целевым объектам:

- EC2-инстансы
- контейнеры
- IP-адреса
- функции Lambda
- виртуальные устройства

Он может распределять трафик приложения с меняющейся нагрузкой в **одной** зоне доступности или между **несколькими** зонами доступности

# Типы балансировщиков ELB

	<b>Application Load Balancer</b>	<b>Network Load Balancer</b>	<b>Gateway Load Balancer</b>	<b>Classic Load Balancer</b>
Уровень	Layer 7	Layer 4	Layer 3, 4	Layer 4/7
Тип цели	IP, инстанс, Lambda	IP, инстанс, ALB	IP, инстанс	
Проверки доступности	HTTP, HTTPS, gRPC	TCP, HTTP, HTTPS	TCP, HTTP, HTTPS	TCP, SSL / TLS, HTTP, HTTPS

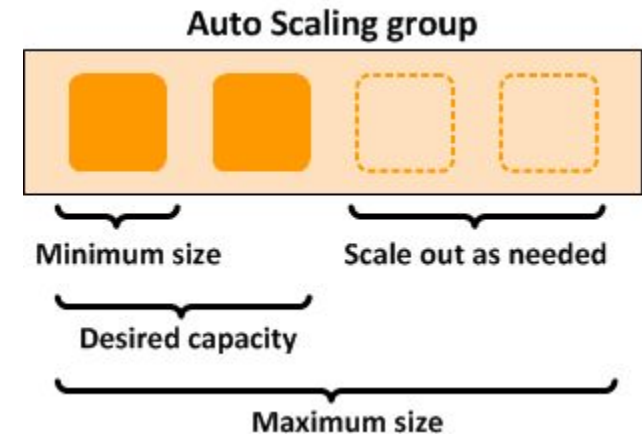


# Auto Scaling Groups

Автоматический запуск новых инстансов Amazon EC2 при повышении спроса и автоматическое отключение ненужных инстансов Amazon EC2 для экономии средств, когда спрос снижается

Состоит из трёх компонентов:

- **Groups** — логическая группа EC2-инстансов одной функциональности
- **Configuration templates** — шаблон конфигурации, применяемый на членов группы (instance type, keypair, SG, AMI ID и т. д.)
- **Scaling options** — политики расширения и сжатия



---

# Scaling Options of Auto Scaling Groups

- **Maintain current instances** — определение количества EC2-инстансов для текущей работы в обычном режиме, определение health-check
- **Ручное масштабирование** — определение количества min и max EC2-инстансов
- **Масштабирование по расписанию**
- **Масштабирование по требованию** — определение параметров, при которых будет произведено расширение ресурсов

## Типы балансировщиков УС

	<b>Application Load Balancer</b>	<b>Network Load Balancer</b>
Уровень	Layer 7	Layer 4
Тип цели	IP, instance	instance
Проверки работоспособности	HTTP, HTTPS	TCP, HTTP

# Типы масштабирования УС

- **Ручное масштабирование** — определение фиксированного количества ВМ
- **Автоматическое масштабирование** — среднее значение метрики не сильно отклонялось от целевого:
  - **региональная**
  - **зональная**

[Автоматически масштабируемая группа ВМ Yandex Cloud](#)



# Итоги



# Итоги

Сегодня мы изучили:

- какие ВМ существуют — EC2-инстансы и Cloud Compute — и их типы
- как организовано хранение, какие виды storage бывают, чем отличаются
- балансировщики нагрузки — какие бывают, что умеют, как настраивать группы масштабирования



# Домашнее задание



## Домашнее задание

Ваше домашнее задание можно посмотреть [по ссылке](#)

- Вопросы по домашней работе задавайте **в чате** учебной группы
- Задачи можно сдавать **по частям**
- Зачёт по домашней работе проставляется после того, как **приняты все задачи**



**Задавайте вопросы и  
пишите отзыв о лекции!**