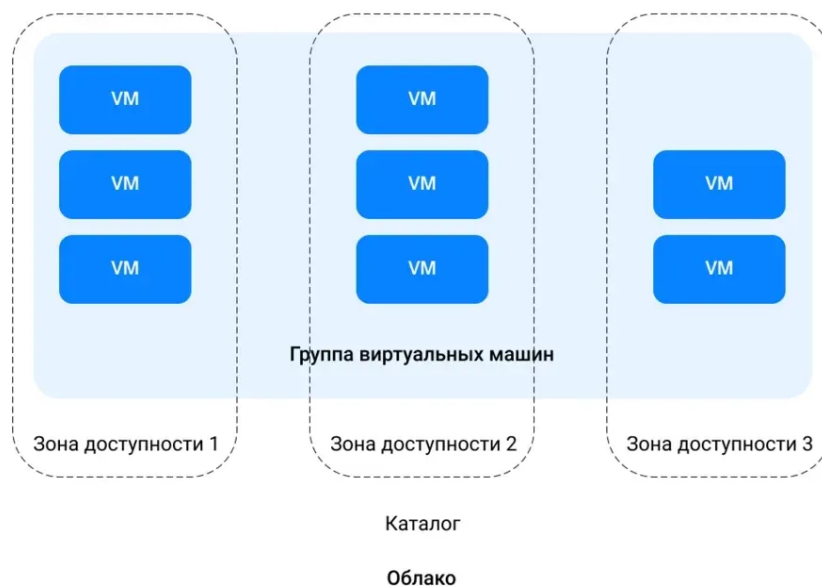


Теория.

Зачем нужны группы виртуальных машин

Управлять большим количеством VM вручную не просто. Всегда есть риск, что вы не отследите программный сбой или пиковую нагрузку, из-за чего сервис — к неудовольствию пользователей — ляжет.

Чтобы избежать таких неприятностей, настройте управление группами VM (или Instance Groups). Сгруппируйте однотипные VM, которые могут находиться в разных зонах доступности, а затем определите, по каким правилам система работает с группами.



Предположим, о вашем стартапе написали в издании Motherboard. На сайт хлынули посетители, нагрузка резко выросла — и VM перестали с ней справляться. В этом случае правильно настроенная система сама создаст достаточно копий VM с приложением. Когда пик интереса пройдёт, система увидит снижение нагрузки и постепенно удалит ненужные копии.

Ещё пример: в одной из VM приложение перестало реагировать на запросы. Система сама определит сбой по заданным правилам и перезапустит эту VM или создаст заново.

Все VM в группе автоматически создаются по шаблону. Вы заполните его параметры при формировании группы. Шаблон описывает конфигурацию машины: какие ей нужны системные ресурсы, как создать дополнительный диск, какие сетевые параметры применить, создавать ли пользователей в системе автоматически и т.д. Создание, обновление и удаление VM в группах выполняется от имени так называемого [сервисного аккаунта](#). Это учетная запись со специфичным набором привилегий (например, административным). Группе VM можно присвоить только один сервисный аккаунт, созданный в том же самом каталоге. Вы также можете использовать сервисный аккаунт для работы с другими API Yandex Cloud (например для интеграции групп VM с сетевым балансировщиком).

Автоматическое восстановление

Ни одно приложение не работает идеально. Например, если сервис из-за программного сбоя начнёт создавать множество временных файлов, на диске рано или поздно закончится свободное место. Работа сервиса прекратится. Пользователи, чьи запросы обслуживает VM, будут видеть сообщение об ошибке. Чтобы VM простаивала как можно меньше, Instance Groups регулярно проверяет состояние VM или отзывчивость приложения. Обнаружив неполадки, сервис действует по выбранному вами сценарию: перезапускает VM или создаёт новую. Способ автоматического восстановления при сбое зависит от того, как вы настроили политику развёртывания:

- Если вы разрешили **превышать** целевой размер группы (поле **Добавлять выше целевого значения**), Instance Groups будет создавать VM вместо не прошедших проверку.
- Если вы разрешили **уменьшать** целевой размер группы (поле **Уменьшать относительно целевого значения**), Instance Groups перезагрузит VM. Иногда для устранения проблемы этого достаточно. Если проблема из примера выше в том, что в папке `/tmp` скопилось много файлов, при перезапуске системы папка автоматически очистится.

Если вы не знаете заранее, достаточно ли перезагрузки VM, комбинируйте оба способа восстановления: используйте сразу два параметра.

Допустим, вы разрешили и превышать, и уменьшать целевой размер группы на одну машину. Когда одна из VM не пройдет проверку, Instance Groups начнет одновременно перезапускать эту машину и создавать новую. VM, которая первая пройдет все проверки, начнет работать, а вторая будет удалена.

Старые машины не удаляются до тех пор, пока не созданы новые. А если в процессе создания новой VM все машины в группе станут работоспособными, то сервис отменит её создание.

Автоматическое восстановление прерываемых VM начнётся только тогда, когда в зоне доступности будет достаточно вычислительных ресурсов. Иногда это занимает немало времени.

Автоматическое масштабирование

Вы разработали и запустили веб-сервис, дали к нему ранний доступ парочке популярных блогеров и со дня на день ждёте наплыва посетителей. Теперь надо сделать так, чтобы сервис продолжил работать при пиковой посещаемости как ни в чём не бывало.

Для этого настройте автоматическое масштабирование группы VM.

Система сама будет отслеживать потребность в VM и добавлять их, а при снижении нагрузки — убирать лишние, чтобы экономить ресурсы и деньги.

Вот как это работает:

1. Создайте группу VM.
2. Укажите, какие метрики системе отслеживать, чтобы вовремя добавлять или убирать VM. Обычно это нагрузка CPU: при загруженном на 100% процессоре сервис попросту перестаёт отзываться на действия посетителей. Вы можете использовать и свои метрики (например время ответа сервиса).
3. Укажите целевое значение метрики. Например **загрузка CPU** — в среднем не больше 50%.

Однако нагрузка на ресурсы бывает неравномерной. Например, ваш сервис мониторинга репутации в соцсетях хранит копии публикаций в микроблогах, на форумах, сайтах с отзывами и т. д. Пользователи часто будут делать выгрузки для отчётности, что подразумевает запросы по очень большому диапазону записей в очень большой базе данных. Поэтому среднее значение метрики может резко меняться. Но если после каждого всплеска и спада нагрузки создавать и удалять ВМ — это тоже приведёт к расходу ресурсов. Поэтому количество ВМ регулируется при помощи нескольких переменных:

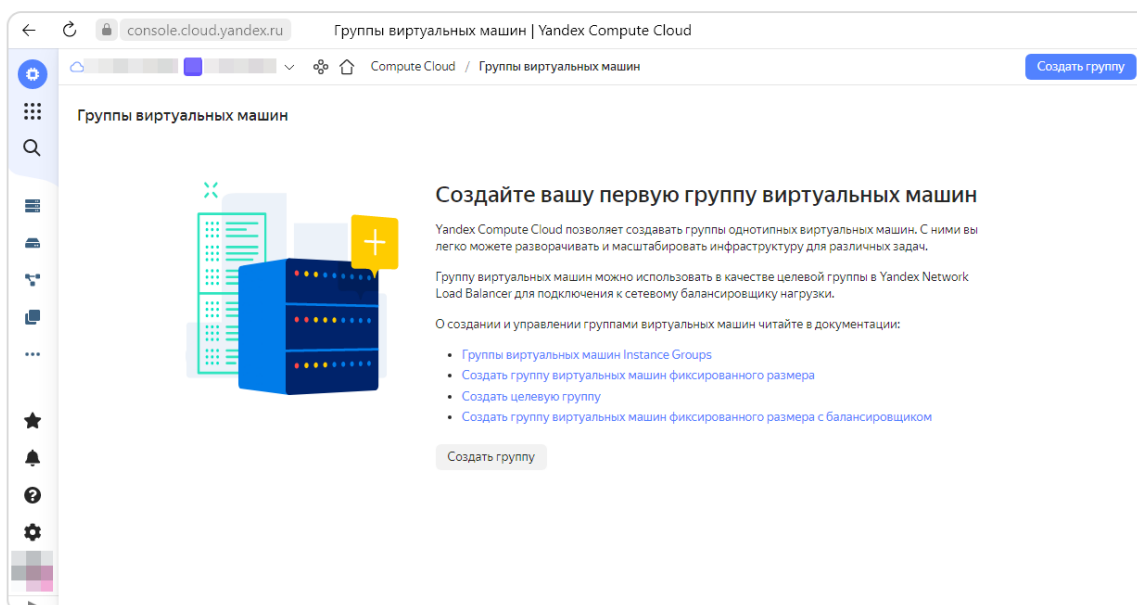
- **Период стабилизации.** После увеличения количества машин в группе ВМ не удаляются сразу и сохраняются заданное время, даже если среднее значение метрики стало заметно ниже её целевого значения. В этом случае при повторном всплеске нагрузки уже будет доступна хотя бы одна дополнительная ВМ, которая перехватит часть запросов.
- **Период прогрева.** Запуск ВМ иногда приводит к аномально высокой нагрузке. Период прогрева позволяет игнорировать ее в течение заданного времени.
- **Период усреднения.** Instance Groups использует усредненные значения всех метрик и игнорирует резкие скачки нагрузки CPU за некоторую единицу времени. Вне периодов стабилизации и прогрева Instance Groups несколько раз в минуту измеряет нагрузку CPU на каждой ВМ и усредняет значения за время, указанное как период усреднения. Затем выполняет дополнительное усреднение по зонам доступности. Если по результатам расчёта нужно создать ещё одну ВМ — Instance Groups запустит этот процесс и начнёт рассчитывать среднюю нагрузку заново.

Вы также можете установить в [Yandex Monitoring](#) пользовательские метрики. Например, среднее время ответа сервиса. Укажите имя метрики и ее целевое значение. Если оно будет превышено, Instance Groups создаст дополнительные машины для распределения нагрузки.

Практическая работа. Создание группы виртуальных машин

Иногда вам требуется не автоматическое масштабирование, а автоматическое восстановление VM. Например, если вы отлаживаете работу веб-сервиса, который периодически падает. Для этого подойдут группы VM фиксированного размера. Давайте создадим и настроим такую группу.

В консоли управления откройте раздел **Compute Cloud**, перейдите на вкладку **Группы виртуальных машин** и нажмите кнопку **Создать группу**.



Откроется страница **Создание группы виртуальных машин**.

В блоке **Базовые параметры** введите имя и описание группы VM. Создайте новый сервисный аккаунт. Чтобы иметь возможность создавать, обновлять и удалять VM в группе, назначьте сервисному аккаунту роль `editor`. По умолчанию все операции в группе VM выполняются от имени сервисного аккаунта.

← console.cloud.yandex.ru Создание группы виртуальных машин | Yandex Compute Cloud

Compute Cloud / Группы виртуальных машин / Создать

Создание группы виртуальных машин

Базовые параметры

Имя ? test-ig-vm

Описание ?

Сервисный аккаунт ? Instance-group или Создать новый

Защита от удаления

Распределение

Зона доступности ?

- ☒ ru-central1-a
- ☐ ru-central1-b
- ☒ ru-central1-c

⚠ Количество свободных ресурсов в зоне ru-central1-c ограничено в связи с плановым выводом зоны из эксплуатации. Рекомендуем выбрать другую зону доступности. [Подробнее о планах](#)

ВМ группы могут находиться в разных зонах и регионах. В блоке **Распределение** выберите две зоны доступности, чтобы обеспечить доступность сервиса, если в одной из них случится сбой. В блоке **Шаблон виртуальной машины** нажмите кнопку **Задать**.

← console.cloud.yandex.ru Создание группы виртуальных машин | Yandex Compute Cloud

Compute Cloud / Группы виртуальных машин / Создать

Шаблон виртуальной машины

Отсутствует шаблон ВМ

Вам необходимо задать конфигурацию базовой виртуальной машины

Задать

В процессе создания и обновления разрешено

Добавлять выше целевого значения ? 0

Уменьшать относительно целевого значения ? 1

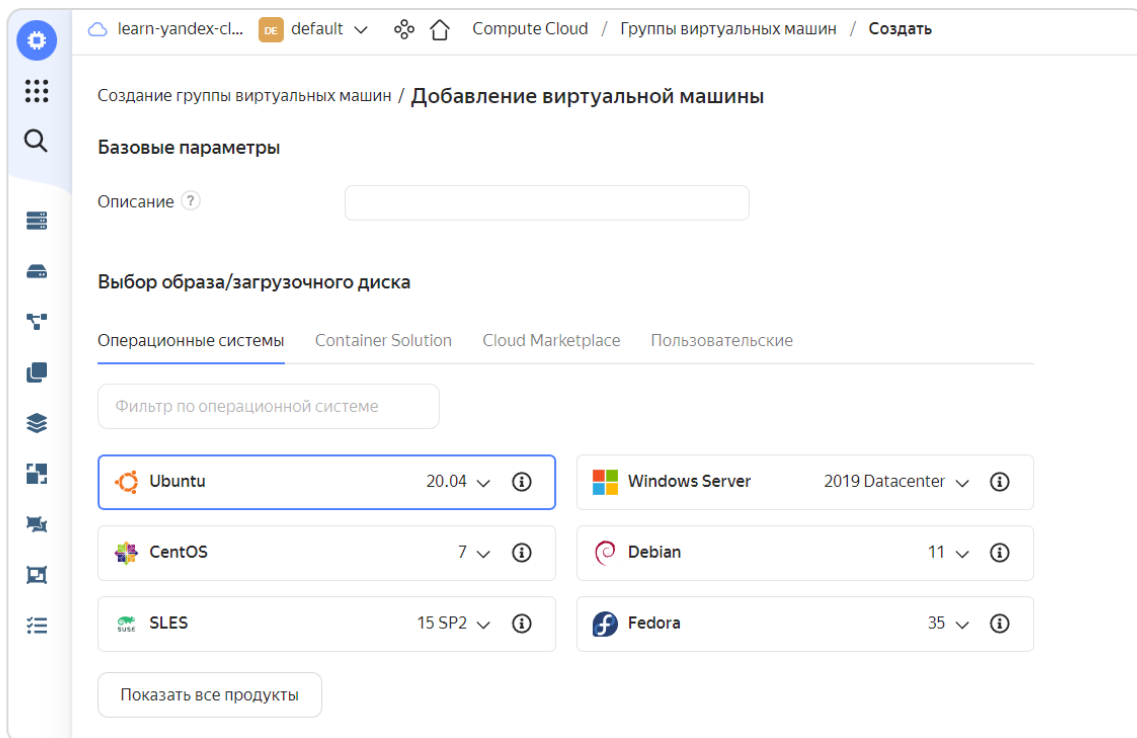
Одновременно создавать ?

Время запуска ? 0 секунд

Одновременно останавливать ?

Останавливать машины по стратегии ? Принудительная Деликатная

Шаблон создается так же, как и сама VM. В блоке **Базовые параметры** введите описание шаблона конфигурации, затем в блоке **Выбор образа/загрузочного диска** на вкладке **Операционные системы** выберите **Ubuntu**.



В блоках **Диски** и **Вычислительные ресурсы** для загрузочного диска оставьте значения по умолчанию. В блоке **Сетевые настройки** выберите существующую сеть и подсеть или создайте новые. В блоке **Доступ** выберите существующий или создайте новый сервисный аккаунт, укажите логин, вставьте в поле **SSH-ключ** содержимое файла с публичным ключом, доступ к серийной консоли не разрешайте.

Сохраните параметры и вы вернётесь на страницу **Создание группы виртуальных машин**.

В блоке **В процессе создания и обновления разрешено** установите политику развертывания:

Добавлять выше целевого значения (на сколько VM можно превышать размер группы) — 2.

Уменьшать относительно целевого значения (на сколько VM можно уменьшать размер группы) — 1.

Одновременно создавать (сколько VM можно сразу создавать в группе) — 2.

Время запуска (сколько времени должно пройти, прежде чем будут пройдены все проверки состояния и VM начнет получать нагрузку) — 2 минуты.

Одновременно останавливать (сколько VM можно сразу удалять) — 1.

Останавливать машины по стратегии — Принудительная. При принудительной стратегии Instance Groups самостоятельно выбирает, какие VM остановить.

console.cloud.yandex.ru Создание группы виртуальных машин | Yandex Compute Cloud

Compute Cloud / Группы виртуальных машин / Создать

В процессе создания и обновления разрешено

Добавлять выше целевого значения ? 2

Уменьшать относительно целевого значения ? 1

Одновременно создавать ? 2

Время запуска ? 2 минут

Одновременно останавливать ? 1

Останавливать машины по стратегии ? Принудительная Деликатная

Масштабирование

Тип ? Фиксированный

Размер 3

Тарифы и цены

Intel Ice Lake. 100% vCPU

Intel Ice Lake. RAM

Windows Server Datacenter. 100% vCPU

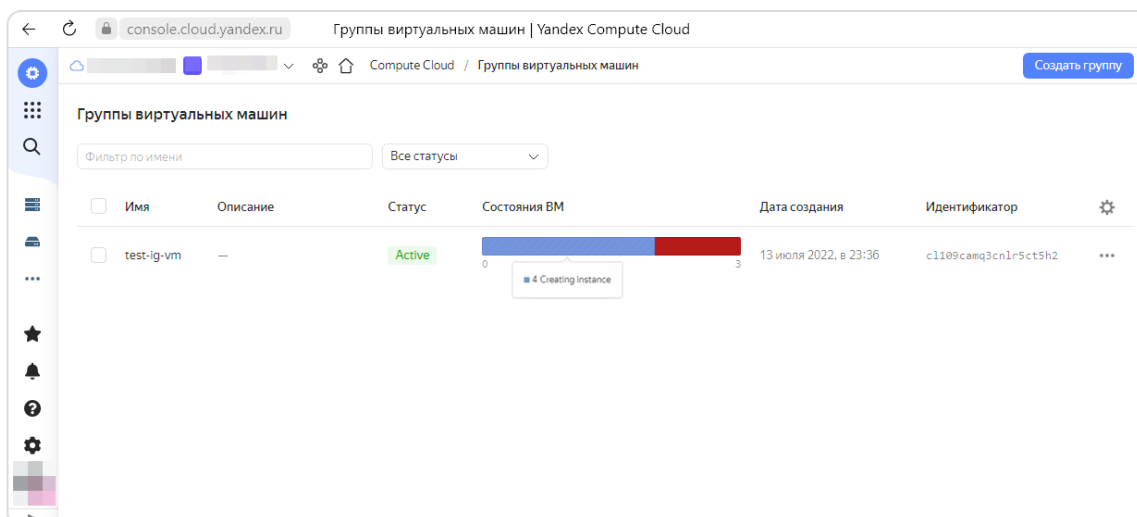
Публичный IP-адрес

Стандартное сетевое хранилище (HDD)

В блоке **Масштабирование** выберите **фиксированный тип**, **Размер** (количество VM) — 3.

В блоке **Интеграция с Load Balancer** оставьте опцию **Создать целевую группу** выключенной. Не включайте пока проверку состояний, которая позволяет Instance Groups получать сведения о состоянии VM.

Нажмите кнопку **Создать** и вернитесь на страницу **Группы виртуальных машин**. В правом нижнем углу появится сообщение «Группа виртуальных машин создаётся». Одновременно можно создавать не более двух VM. Поэтому сначала будут созданы две VM, потом — третья.

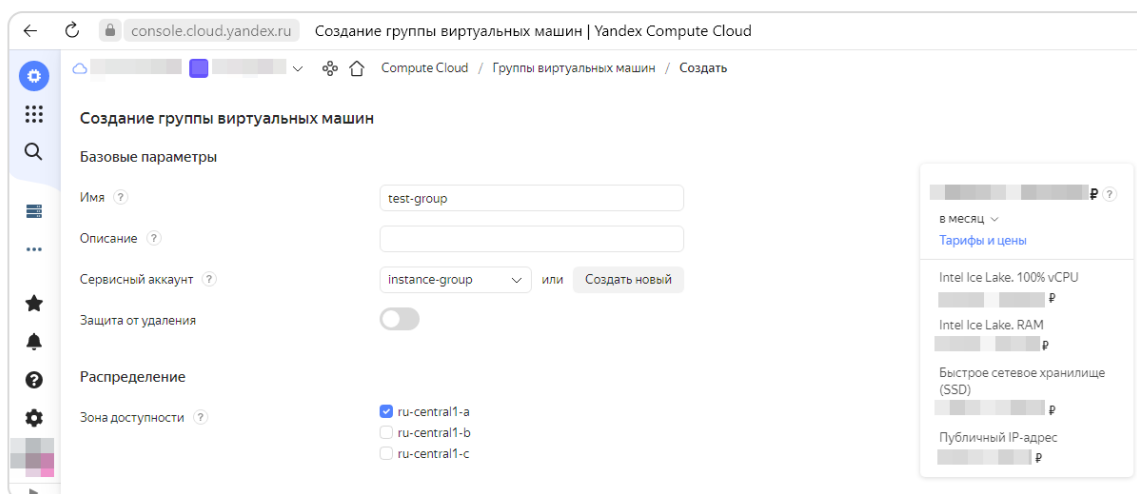


После того как вы создали группу, протестируйте включение и выключение всех машин сразу. Обратите внимание: в соответствии с настройками сервис инициирует запуск не более двух машин одновременно. Третья VM будет оставаться остановленной. Как только первая будет запущена, один слот на запуск освободится, поэтому сразу будет инициирован запуск третьей и последней VM.

Практическая работа. Автоматическое масштабирование под нагрузкой

Давайте разберёмся, как обеспечить доступность сервиса под высокой нагрузкой. Вы уже научились создавать группы VM. Теперь создадим автоматически масштабируемую группу VM.

В консоли управления откройте раздел **Compute Cloud**. Перейдите на вкладку **Группы виртуальных машин** и нажмите кнопку **Создать группу**. Задайте имя группе VM.



Создайте [сервисный аккаунт](#). Чтобы иметь возможность создавать, обновлять и удалять VM в группе, назначьте сервисному аккаунту роль `editor`. По умолчанию все операции в группе VM выполняются от имени сервисного аккаунта.

В блоке **Распределение** выберите только одну зону доступности.

В блоке **Шаблон виртуальной машины** нажмите кнопку **Задать**. В открывшемся окне выберите:

- ОС: Ubuntu 20.04.

- Размер загрузочного диска: 50 ГБ.
- Тип загрузочного диска: SSD.

Остальные параметры — по умолчанию. Не забудьте добавить публичный SSH-ключ. Он понадобится нам на следующем практическом занятии.

Создание группы виртуальных машин | Yandex Compute Cloud

Шаблон виртуальной машины

Вычислительные ресурсы

Платформа..... Intel Ice Lake

Гарантированная доля vCPU..... 100%

vCPU..... 2

RAM..... 2 ГБ

Сетевые настройки

Сеть..... intro2

Подсети (ru-central1-a)..... intro2-ru-central1-a

Публичный IP-адрес..... Автоматически

Группы безопасности..... free

Диски

Ubuntu 20.04 LTS, 50 ГБ, SSD **Загрузочный**

Доступ

Доступ к серийной консоли..... **запрещен**

В процессе создания и обновления разрешено

Добавлять выше целевого значения ? 0

Уменьшать относительно целевого значения ? 1

Одновременно создавать ?

Время запуска ? 0 секунд

Одновременно останавливать ?

Останавливать машины по стратегии ? Принудительная Деликатная

в месяц

Тарифы и цены

Intel Ice Lake, 100% vCPU

Intel Ice Lake, RAM

Быстрое сетевое хранилище (SSD)

Публичный IP-адрес

В блоке **В процессе создания и обновления разрешено** оставьте параметры по умолчанию.

Перейдите к блоку **Масштабирование** и выберите тип **Автоматический**.

Создание группы виртуальных машин | Yandex Compute Cloud

Масштабирование

Тип ? Автоматический

Автоматическое масштабирование групп работает в экспериментальном режиме

Тип автомасштабирования ? Зональное Региональное

Минимальное количество VM в зоне ? 2

Максимальный размер группы ? 4

Промежуток измерения нагрузки ? 60 секунд

Время на разогрев VM ? 3 минут

Период стабилизации ? 5 минут

Начальный размер группы ? 4

в месяц

Тарифы и цены

Intel Ice Lake, 100% vCPU

Intel Ice Lake, RAM

Быстрое сетевое хранилище (SSD)

Публичный IP-адрес

Задайте параметры масштабирования:

- **Тип автомасштабирования** — зональное. При зональном автомасштабировании количество VM регулируется отдельно в каждой зоне доступности, указанной в настройках группы.
- **Минимальное количество VM в зоне** — 2. Сервис Instance Groups не будет удалять VM в зоне доступности, если их там всего две.
- **Максимальный размер группы** — 4. Instance Groups не будет создавать VM, если их уже четыре. В этот раз размер загрузочного диска VM — 50 ГБ, поэтому с учётом квот на суммарный объём SSD-дисков в одном облаке смогут запуститься четыре VM.
- **Промежуток измерения загрузки** (это период усреднения: время, за которое следует усреднять замеры нагрузки для каждой VM в группе) — 60 секунд.
- **Время на разогрев VM** — 3 минуты. В течение этого времени VM не учитывается в измерении средней нагрузки

на группу. Фактически данное время мы можем определить, измерив, как быстро запускается VM.

- **Период стабилизации** — 5 минут. Отсчитывается с момента, когда Compute Cloud принял последнее решение о том, что количество VM в группе нужно увеличить.
- **Начальный размер группы** — 4. Это количество VM, которое следует создать вместе с группой.

В блоке **Метрики** укажите:

- **Тип** — CPU.

Целевой уровень загрузки CPU, % — 80. Instance Groups будет управлять размером группы так, чтобы поддерживать указанную нагрузку CPU.

console.cloud.yandex.ru Создание группы виртуальных машин | Yandex Compute Cloud

Compute Cloud / Группы виртуальных машин / Создать

Метрики

CPU

Метрика: CPU

Целевой уровень загрузки CPU, %: 80

Добавить метрику

Интеграция с Network Load Balancer

Создать целевую группу: ☐

Интеграция с Application Load Balancer

Создать целевую группу: ☐

Проверка состояний

Активировать: ☐

Пользовательские переменные

Ключ: Значение

Добавить поле

Создать

в месяц

Тарифы и цены

Intel Ice Lake, 100% vCPU

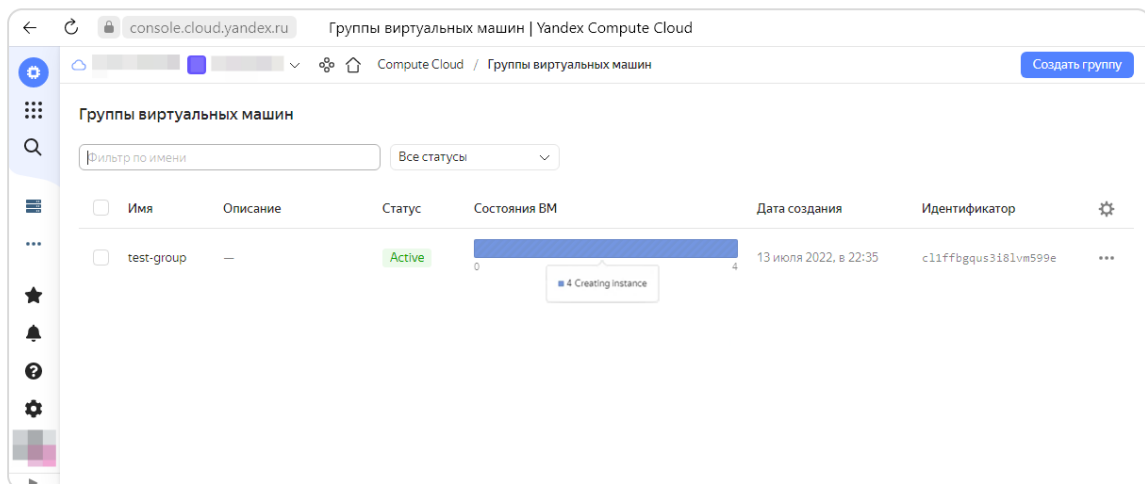
Intel Ice Lake, RAM

Быстрое сетевое хранилище (SSD)

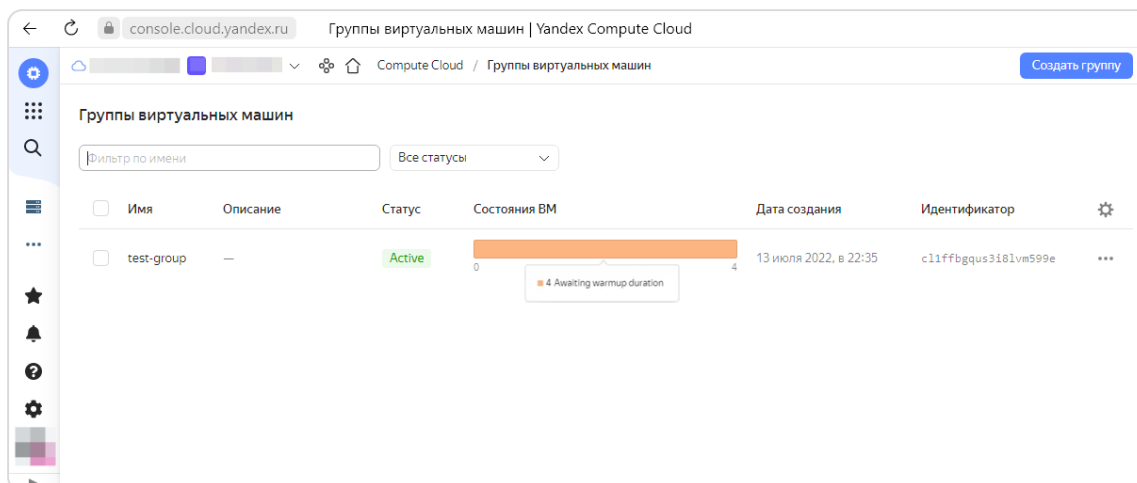
Публичный IP-адрес

Нажмите кнопку **Создать**. Сервис начнет создавать ВМ. После создания статус группы изменится на **Active**. Обратите внимание, как меняются **Состояния ВМ**.

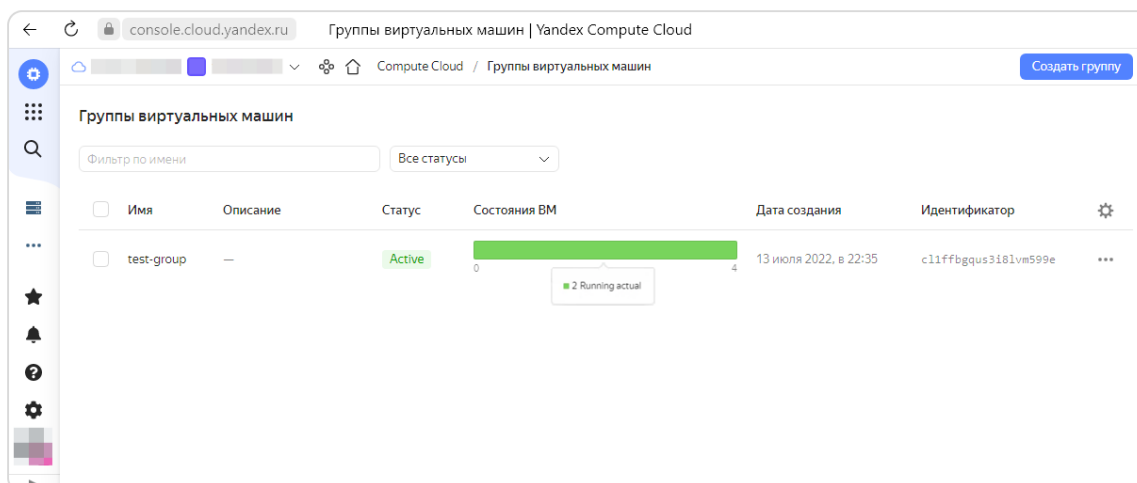
Creating instance — ВМ создаётся и запускается.



Awaiting warmup duration — ВМ начинает принимать сетевой трафик. В этом статусе ВМ находится в течение периода прогрева, указанного в настройках автоматического масштабирования. Значения метрик ВМ в этом статусе заменяются средними значениями ВМ из той же зоны доступности.



Running actual — ВМ запущена, на неё подается сетевой трафик, пользовательские приложения работают.



Группа ВМ готова принимать рабочую нагрузку.

Практическая работа. Воссоздание виртуальных машин в группе

Давайте симитируем рост нагрузки на VM и посмотрим, как сервис на это отреагирует.

1. В консоли управления откройте раздел **Compute Cloud** и перейдите на страницу группы VM, которую вы создали на прошлом практическом занятии. Откройте в двух вкладках браузера страницы **Группы виртуальных машин** и **Мониторинг** (открывается из раздела **Группы виртуальных машин**).
2. После запуска группы зайдите по SSH на каждую из двух VM и установите приложение для стресс-тестирования Linux-систем. Для этого выполните команду:

Скопировать код

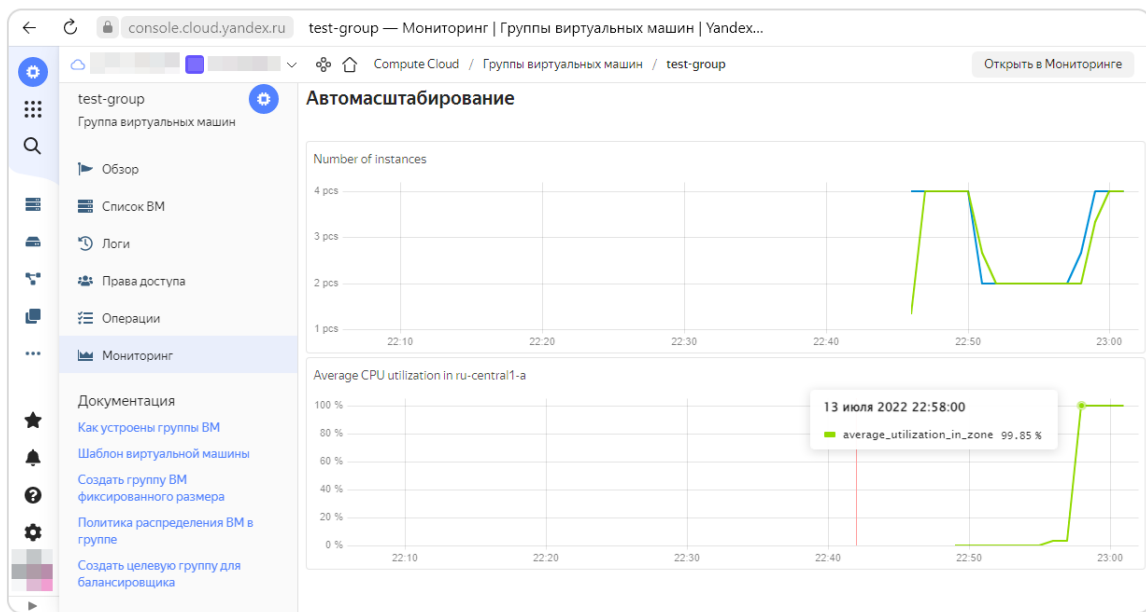
```
sudo apt-get install stress
```

3. После этого для каждой VM запустите установленное приложение:

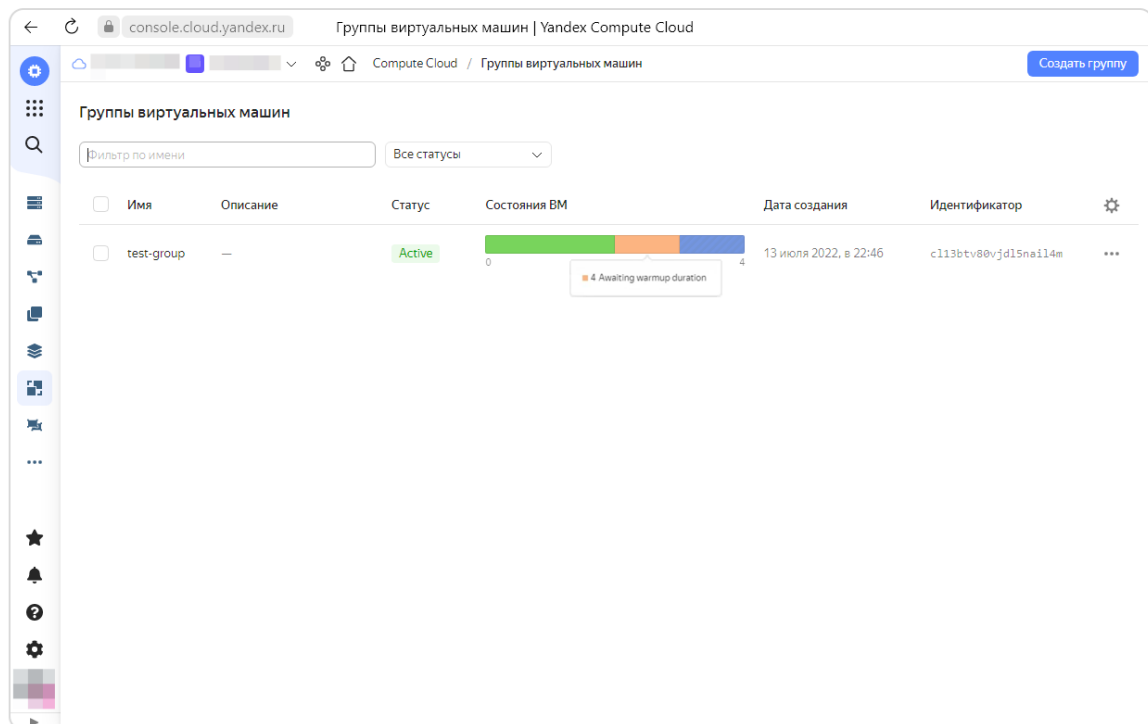
Скопировать кодBASH

```
stress -c 2
```

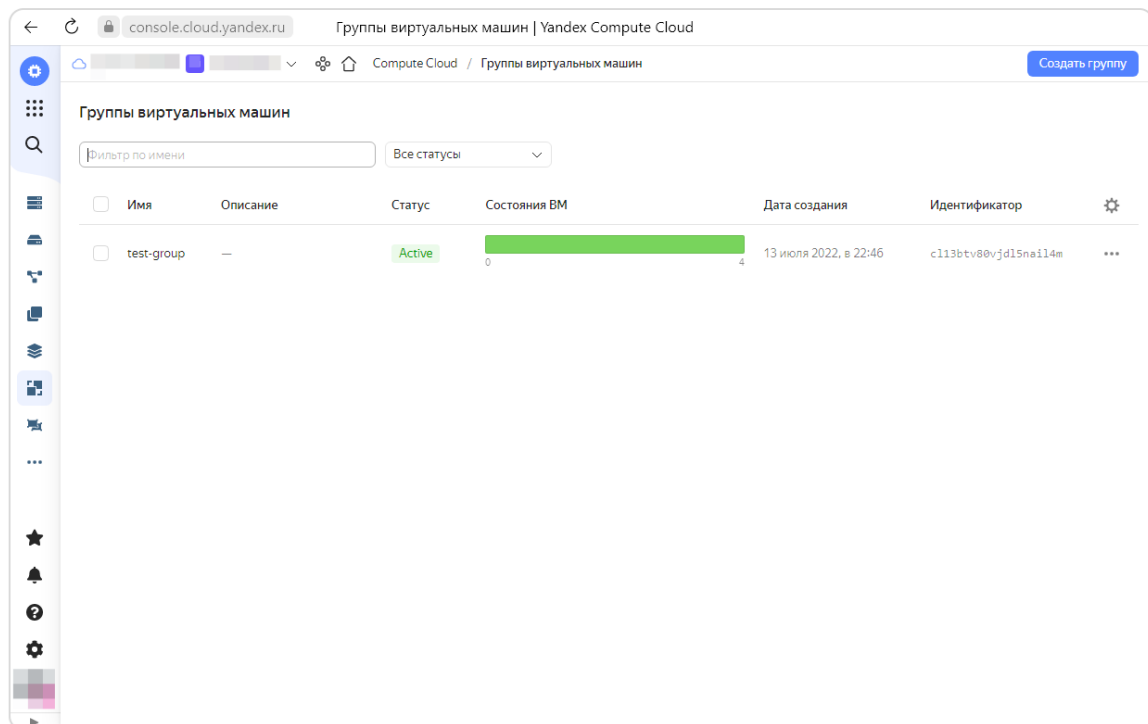
Аргумент `-c` означает, что при тестировании будет нагружаться процессор, а число после аргумента задаёт количество ядер процессора, которые будут нагружаться. Чтобы эксперимент удался — укажите количество ядер, которое вы выбрали в шаблоне VM. На вкладке со страницей мониторинга на графике **Average CPU utilization in ru-central1-a** следите за тем, как усреднённое значение нагрузки будет постепенно расти.



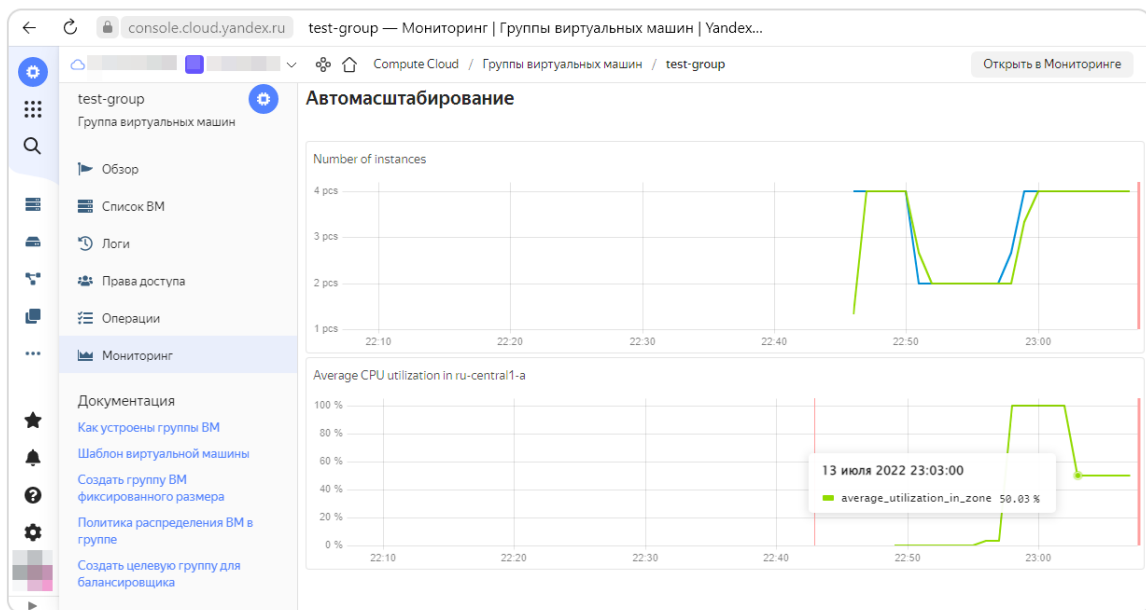
Как только усреднённое значение нагрузки превысит порог, сервис Instance Groups начнёт прогревать две дополнительные VM и вводить их в строй. Это будет видно на странице **Группы виртуальных машин**.



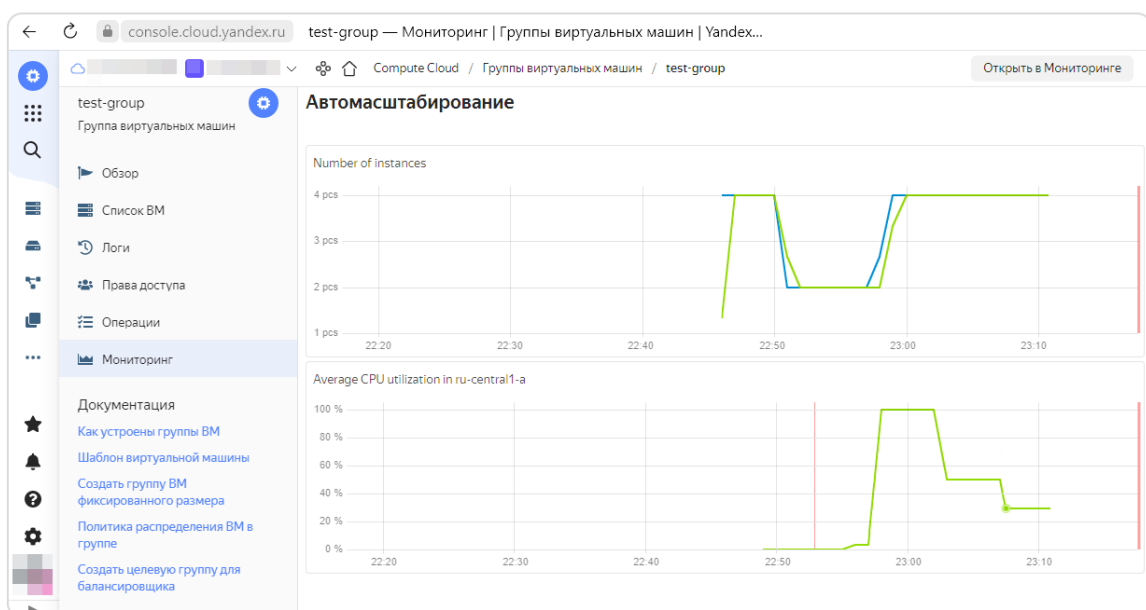
Поскольку стресс-тест не остановлен, сервис завершает запуск двух VM.



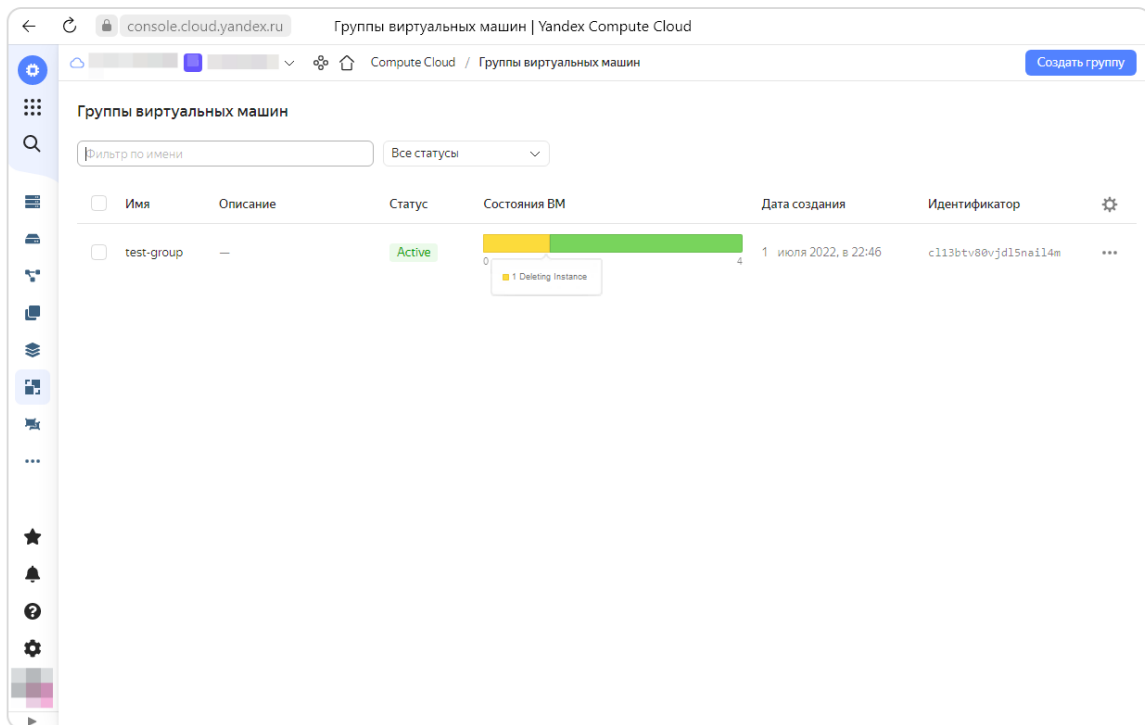
Через некоторое время усреднённое значение нагрузки процессоров в группе упадёт до 50%, поскольку первая половина ВМ загружена полностью, а вторая не загружена вовсе.



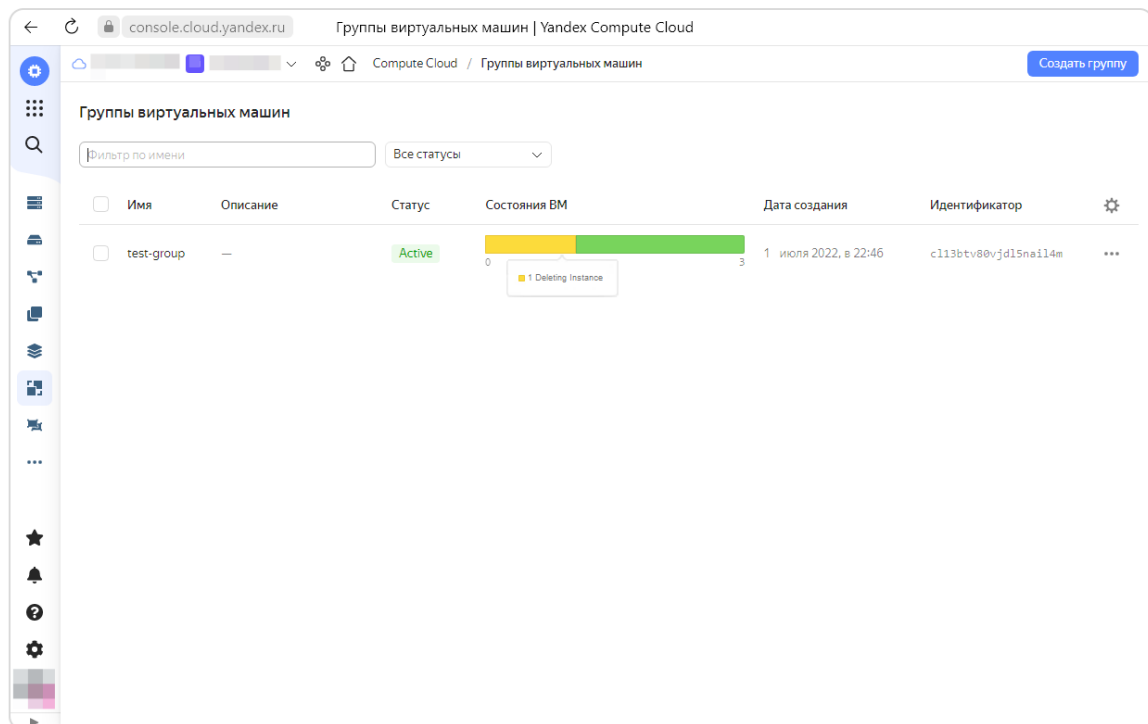
Остановите работу стресс-теста на первой VM. В командной строке используйте сочетание клавиш **Ctrl + C**.



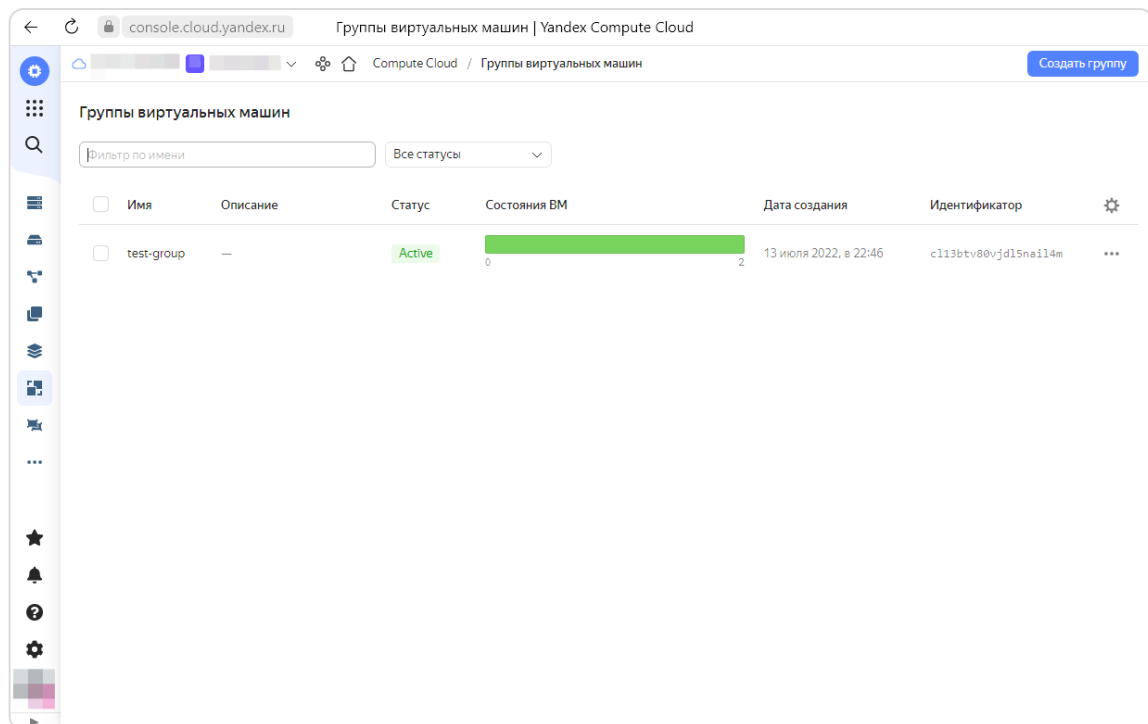
Через некоторое время усреднённое значение достигнет 25%, тогда Instance Groups удалит лишнюю VM:



Остановите второй стресс-тест. Через некоторое время после того, как усредненное значение достигнет нуля, Instance Groups удалит вторую дополнительную VM.



При минимальной нагрузке остаются работать две машины:



Вот так при растущей нагрузке группа VM автоматически масштабируется, чтобы обеспечить доступность ресурса.