

Масштабируемость и отказоустойчивость Docker-кластеров

Балашов Антон



План лекции

- Масштабируемость и отказоусточивость приложений
- Кластеры Docker
- Балансировка нагрузки приложений в кластере

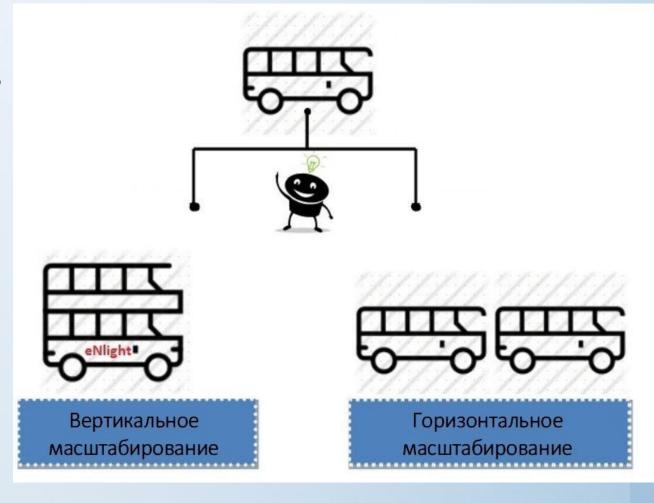
Масштабируемость

Масштабируемость при облачных вычислениях — это возможность быстро и без труда увеличить или уменьшить размер либо мощность ИТ-решения или ресурса. В то время как термин «масштабируемость» может означать способность любой системы справиться с растущим объемом работы, в контексте горизонтального и вертикального масштабирования речь часто идет о базах данных и больших объемах данных.

Масштабируемость позволяет адаптироваться к огромным объемам разнообразных данных и управлять ими, изменяя объемы данных и шаблоны рабочих нагрузок, которые создаются в облаке, на мобильных устройствах, в социальных сетях и источниках больших данных.

Масштабируемость

- Вертикальное масштабирование: наращивание или снижение вычислительной мощности или ресурсов по мере необходимости. Для автоматической адаптации к требованиям рабочих нагрузок применяется либо изменение уровней производительности, либо эластичные пулы ресурсов.
- Горизонтальное масштабирование: добавление дополнительных мощностей или разделение одного узла на узлы меньшего размера с использованием сегментирования. Обеспечивает ускоренное и более удобное управление приложением на серверах.



Вертикальное масштабирование

Вертикальное масштабирование используется, когда необходимо быстро реагировать на проблемы с производительностью, которые нельзя решить с помощью классической методики оптимизации кода. Вертикальное масштабирование помогает справиться с пиками в рабочих нагрузках, когда текущий уровень производительности не может удовлетворить все требования. Вертикальное увеличение масштаба позволяет добавлять дополнительные ресурсы, чтобы легко адаптироваться к пиковым рабочим нагрузкам. Затем, если ресурсы больше не нужны, можно выполнить вертикальное уменьшение масштаба, чтобы вернуться к исходному состоянию и сократить затраты.

Вертикальное масштабирование выполняется в следующих случаях:

- Рабочие нагрузки достигают определенного ограничения производительности, например ограничения, касающегося ЦП или операций ввода-вывода.
- Нужно быстро реагировать, чтобы устранить проблемы с производительностью, которые нельзя решить путем классической оптимизации базы данных.
- Требуется решение, позволяющее изменять уровни служб, чтобы адаптироваться к изменению требований к задержке.

Горизонтальное масштабирование

Горизонтальное масштабирование применяют когда не удается получить достаточно ресурсов для рабочих нагрузок даже на самых высоких уровнях производительности. При горизонтальном масштабировании копии приложения распределяются между серверами. Масштаб каждого сервера можно вертикально увеличивать или уменьшать по отдельности.

Горизонтальное масштабирование выполняется в следующих случаях:

- Географически распределенные приложения, каждое из которых должно работать в своем регионе.
- Глобальный сценарий сегментирования (например, балансировка нагрузки) с большим количеством географически распределенных клиентов.
- Если ограничения производительности превышаются даже на самых высоких уровнях производительности службы.

Отказоустойчивость приложений

Отказоустойчивость — свойство технической системы сохранять свою работоспособность после отказа одной или нескольких её составных частей. Отказоустойчивость определяется количеством единичных отказов составных частей (элементов) системы, после наступления которых сохраняется работоспособность системы в целом. Базовый уровень отказоустойчивости подразумевает защиту от отказа одного любого элемента. Поэтому основной способ повышения отказоустойчивости это избыточность.

Реализация отказоустойчивости:

• DR (disaster recovery)- копия системы может быть поднята в другом ЦОД. Существует ряд требований к реализации такого сценария.

Программно реализованная отказоустойчивость для каждого из компонент и их взаимодействия

Программно реализованная отказоустойчивость для каждого из компонент и их взаимодействия

- Low level fault tolerant services система должна состоять из независимых подсистем, каждая из которых должна быть отказоустойчивой.
- Single point of failure Избежание архитектуры, где вся система рушиться при падении одного из компонент. Достичь это можно либо используя принцип избыточности, либо делать компоненты по возможности независимыми, чтобы при падении одного из компонент, переставала работать только часть функциональности, а остальные части системы продолжали работать.
- **Избыточность (Redundancy) в** системе существует несколько копий компонент. И при падении одного из них система продолжает работать. При таком подходе проектирования можно выделить следующие стратегии:

Стратегии избыточности

- Active-Active- система одновременно работает с двумя идентичными компонентами. Запрос на обработку поступает одновременно на два компонента, одновременно обрабатывающие запрос. Если один из таких компонентов рушится, то вся работа автоматически переходит только на второй компонент. В качестве минусов можно выделить, увеличенное количество трафика и время на его обработку, а так же дополнительная серверная инфраструктура в постоянно рабочем режиме потребления ресурсов.
- Active-Passive- только один постоянно работающий компонент, в случае падения автоматически поднимается второй, восстанавливает состояние и берет на себя всю работу. Минусом является время восстановления состояния

• Балансировка нагрузки (Load Balancing) - несколько идентичных компонентов, между которыми распределяется нагрузка по заданному правилу. В отличие от вышеописанной active-active стратегии, здесь каждую задачу выполняет только один компонент. Данный механизм идеально подходит для stateless компонент, иначе для отказоустойчивости вам постоянно придется синхронизировать состояние. Например в случае веб-серверов — делать репликацию сессии. В данном решении очень важно иметь как минимум N+1 redundancy, т.е. если для пиковых нагрузок вам необходимо N работающих на всю катушку компонент, то в вашей системе должно присутствовать N+1 таких компонент, так как иначе, если у вас упадет один из элементов

ЦОД. Классификация Uptime Institute

Известны четыре уровня стандарта Uptime Institute:

- Tier I базовая инфраструктура без резервирования;
- Tier II инфраструктура с резервными мощностями;
- Tier III инфраструктура, поддерживающая параллельный ремонт;
- Tier IV отказоустойчивая инфраструктура.

Каждый следующий уровень сертификации включает в себя требования для всех предыдущих уровней.

ЦОД. Классификация Uptime Institute

Дата-центр первого уровня предоставляет выделенную инфраструктуру для поддержки IT-систем за пределами офиса, источник бесперебойного питания для фильтрации скачков напряжения и обработки кратковременных отключений, специальное охлаждающее оборудование, которое продолжает функционировать и по завершении рабочего дня, генератор для защиты систем от продолжительных отключений электроэнергии.

ЦОДы второго уровня включают в себя резервные возможности для критически важных компонентов в целях обеспечения ремонта и повышенной защиты IT-процессов от сбоев. Резервируемые системы включают в себя оборудование для питания и охлаждения, такие как источники бесперебойного питания, чиллеры или насосы, а также генераторы.

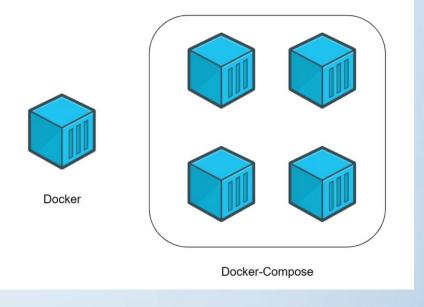
Центр обработки данных Tier III не требует прекращения работы оборудования для замены «железа» и обслуживания. К компонентам добавляется резервный канал питания и охлаждения так, чтобы каждый элемент, необходимый для поддержки IT-системы, можно было отключить, и это не сказалось на работе.

ЦОДы Tier IV в дополнение ко всем особенностям предыдущего уровня характеризуются повышенной (еще большей, чем у Tier III) отказоустойчивостью, то есть сбои отдельных элементов или перебои резервного канала не сказываются на IT-операциях.

Принято считать, что ожидаемый уровень безотказной работы дата-центра Tier I составляет 99,671% (1729 минут годового простоя); Tier 2 — 99,741% (1361 минут годового простоя); Tier III — 99,982% (95 минут годового простоя); Tier IV — 99,995% (26 минут годового простоя).

Кластеры Docker

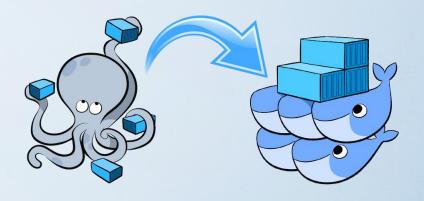
Docker Compose — это инструментальное средство, входяц состав Docker. Оно предназначено для решения задач, связ с развёртыванием проектов, состоящих из нескольких контейнеров.



Разница между Docker и Docker Compose

- Docker применяется для управления отдельными контейнерами (сервисами), из которых состоит приложение.
- Docker Compose используется для одновременного управления несколькими контейнерами, входящими в состав приложения. Этот инструмент предлагает те же возможности, что и Docker, но позволяет работать с более сложными приложениями.

Возможности Docker Compose



Несколько изолированных сред на одном хосте

- Возможно создавать несколько сред окружения на одном хосте, используя название проекта в различных контекстах:
- Создать на одном хосте несколько копий одного и того же окружения.
- Изолировать разные проекты на хосте, которые могут использовать сервисы с одинаковыми названиями.

Пересоздаются только измененные контейнеры

Сотрозе кэширует конфигурацию, которая была использована для создания контейнера. Если перезапустить контейнер без изменений, то будут использованы существующие файлы. Такое повторное использование обеспечивает быстрое внесение изменений в проект.

Docker Compose установка

- Windows:
 - Входит в пакет Docker Desktop
- Linux, MacOs:
 - https://docs.docker.com/compose/install/

Возможности Docker Compose

Команды запуска

- Создание контейнера для служб и запуск среды в фоновом режиме docker-compose up -d
- проверить активность среды: docker-compose ps
- Посмотреть логи:
 docker-compose logs
- Приостановить работу контейнерной среды без изменения текущего состояния контейнеров:

docker-compose pause

- Возобновить работу после паузы: docker-compose unpause
- Остановить выполнение контейнера без удаления связанных с ним данных: docker-compose stop
- Удалить контейнеры, сети и тома, которые связаны с контейнерной средой: docker-compose down
- Удалить образ, из которого собирается среда: docker image rm name:tag

Описание среды docker-compose

Файл docker-compose.yml должен начинаться с тега версии.

version: "3.2"

Следует учитывать, что docker-composes работает с сервисами. 1 сервис = 1 контейнер.

Сервисом может быть клиент, сервер, сервер баз данных...

Раздел, в котором будут описаны сервисы, начинается с 'services'.

services:

Название сервиса

whale:

Docker образ, который мы хотим запустить в контейнере

image: docker/whalesay

Команда, которую нужно запустить после скачивания образа.

```
command: ["cowsay", "hello!"]
```

YAML

Основные цели:

- быть понятным человеку;
- поддерживать структуры данных, родные для языков программирования;
- быть переносимым между языками программирования;
- использовать цельную модель данных для поддержки обычного инструментария;
- поддерживать потоковую обработку;
- быть выразительным и расширяемым;
- быть лёгким в реализации и использовании.



УАМЬ (YAML AIN'T MARKUP LANGUAGE — YAML HE ЯВЛЯЕТСЯ ЯЗЫКОМ РАЗМЕТКИ) — ЭТО УДОБНЫЙ ФОРМАТ СЕРИАЛИЗАЦИИ ДАННЫХ ДЛЯ ВСЕХ ЯЗЫКОВ ПРОГРАММИРОВАНИЯ

yaml.org

YAML Cинтаксис

Отступы

• В YAML для разделения информации очень важны отступы. Нужно помнить, что используются только пробелы, табуляция не допускаются.

При отсутствии отступа перед первым объявлением YAML поймет, что это корень (уровень 0) вашего файла.

Ключ/Значение

• Как и в JSON/JS, в YAML есть синтаксис ключ/значение, значением может быть как число, так и строка, так и список

Списки

- Синтаксис JSON: массив строк

people: ['Anne', 'John', 'Max']

- Синтаксис дефиса

people:

- Anne
- John
- Max

Пример среды docker-compose

```
root@LAPTOP-1UADVU1A:~/work/DevOps/Ла63/python-ui-database# tree -a .

- .env
- app
- Dockerfile
- app.py
- requirements.txt
- db
- init.sql
- docker-compose.yml
```

Способы конфигурирования среды

• Переменные окружения

Web:

image: «webapp:\${TAG}»

• .env файл в директории с docekr-compose.yml

cat.env

MYSQL_USER=root

MYSQL_ROOT_PASSWORD=root

MYSQL_HOST=db

MYSQL_PORT=3306

Способы конфигурирования среды

```
web:
 environment:
   PMA_HOST: db
   PMA_PORT: 3306
   PMA_USER: ${MYSQL_USER}
   PMA_PASSWORD: ${MYSQL_ROOT_PASSWORD}
    env_file - тэг
web:
 env_file:
   - web-variables.env
    env_file - ключ
docker compose —env-file ./config/.env.dev up
```

Балансировка нагрузки приложений в кластере

Балансировка нагрузки подразумевает эффективное распределение входящего сетевого трафика между группой бэкенд-серверов. Задача же регулятора — распределить нагрузку между несколькими установленными бэкенд-серверами.

Балансировка нагрузки помогает масштабировать приложение, справляясь со скачками трафика без увеличения расходов на облако. Она также помогает устранить проблему единой точки отказа. Поскольку нагрузка является распределённой, то в случае сбоя одного из серверов — сервис всё равно продолжит работу.

- Nginx это высокопроизводительный веб-сервер, который также может использоваться в качестве регулятора нагрузки это процесс распределения веб-трафика между несколькими серверами с помощью Nginx. Он гарантирует, что ни один сервер не будет перегружен и что все запросы будут обработаны своевременно. Nginx использует различные алгоритмы для определения оптимального распределения трафика, а также может быть настроен для обеспечения отказоустойчивости в случае выхода из строя одного из серверов.
- НАРгоху серверное программное обеспечение для обеспечения высокой доступности и балансировки нагрузки для ТСР и НТТР приложений, методом распределения входящих запросов на несколько серверов

Nginx

Nginx (NGINX, Engine-X, «Энжин-кс») — это бесплатный веб- и почтовый прокси-сервер с непоточной (асинхронной) архитектурой и открытым кодом.

Разработку Nginx начал в 2002 году Игорь Сысоев для Rambler. А в 2004 году он стал доступен широкому кругу пользователей. С 2011 года серверное ПО начала выпускать уже собственная фирма Игоря, которая спустя 2 года запустила расширенную платную версию продукта (Nginx Plus). Весной 2019 года Nginx была выкуплена крупным американским девелопером F5 Networks.

В отличие от обычного веб-сервера, Nginx не создаёт один поток под каждый запрос, а разделяет его на меньшие однотипные структуры, называемые рабочими соединениями. Каждое такое соединение обрабатывается отдельным рабочим процессом, а после выполнения они сливаются в единый блок, возвращающий результат в основной процесс обработки данных. Одно рабочее соединение может обрабатывать до 1024 запросов одного вида одновременно.

Nginx.Практическое применение

- Отдельный порт/IP. При наличии большого количества статичного контента или файлов для загрузки, можно настроить на отдельном порту или IP-адресу, чтобы осуществлять раздачу. При большом количестве запросов рекомендуется ставить отдельный сервер и подключать к нему Nginx.
- **Акселерированное проксирование**. В таком случае все пользовательские запросы на статичный контент (картинки, простой HTML, JavaScript, CSS-файлы) поступают сначала на Nginx, а он их обрабатывает самостоятельно. При этом никаких изменений исходного кода не требуется.
- **Nginx и FastCGI**. Если поддерживается технология FastCGI, Apache вообще можно не использовать. Но в таком случае может потребоваться модификация кодов скриптов.
- Интерфейс **FastCGI** клиент-серверный протокол взаимодействия веб-сервера и приложения, дальнейшее развитие технологии CGI. По сравнению с CGI является более производительным и безопасным.
- FastCGI снимает множество ограничений CGI-программ. Недостаток CGI-программ в том, что они должны быть перезапущены веб-сервером при каждом запросе, что приводит к понижению производительности. FastCGI, вместо того чтобы создавать новые процессы для каждого нового запроса, использует постоянно запущенные процессы для обработки множества запросов. Это позволяет экономить время.

Round Robin

метод балансировки нагрузки, при котором каждому серверу в кластере предоставляется равная возможность обрабатывать запросы. Этот метод часто используется в веб-серверах, где каждый запрос сервера равномерно распределяется между серверами.

Мощность распределяется поочерёдно, что означает, что у каждого сервера будет своё время для выполнения запроса. Например, если у вас есть три вышестоящих сервера, А, В и С, то балансировщик нагрузки сначала распределит нагрузку на А, затем на В и, наконец, на С, прежде чем перераспределить нагрузку на А. Этот метод довольно прост, но имеет некоторую долю ограничений.

Одно из ограничений заключается в том, что некоторые серверы будут простаивать просто потому, что они будут ждать своей очереди. В нашем примере, если А получит задание и выполнит его за секунду, это будет означать, что он будет простаивать до следующего задания. По умолчанию для распределения нагрузки между серверами Nginx использует именно метод round robin.

Round Robin с добавлением веса

Чтобы решить проблему простоя серверов, мы можем использовать server weights (серверные веса), чтобы указать Nginx, какие серверы должны иметь наибольший приоритет. Weighted Round Robin — один из самых популярных методов балансировки нагрузки, используемых сегодня.

Этот метод предполагает присвоение веса каждому серверу, а затем распределение трафика между серверами на основе этих весов. Это гарантирует, что серверы с большей пропускной способностью получат больше трафика, и поможет предотвратить перегрузку какого-либо из серверов.

Этот метод часто используется в сочетании с другими методами, такими как Session Persistence, для обеспечения равномерного распределения мощностей на все серверы. Сервер приложения с наибольшим параметром веса будет иметь приоритет (больше трафика) по сравнению с сервером с наименьшим числом (весом).

```
upstream app{
   server 10.2.0.100 weight=5;
   server 10.2.0.101 weight=3;
   server 10.2.0.102 weight=1;
}
```

Least Connection

Метод наименьшего числа соединений (Least Connection) — это популярная техника, используемая для равномерного распределения рабочей нагрузки между несколькими серверами. Метод работает путём маршрутизации каждого нового запроса на соединение — на сервер с наименьшим количеством активных соединений. Это гарантирует, что все серверы используются одинаково и ни один из них не перегружен.

```
upstream app{
    least_conn;
    server 10.2.0.100;
    server 10.2.0.101;
    server 10.2.0.102;
}
```

Least Connection с добавлением веса

Данный метод используется для распределения рабочей нагрузки между несколькими вычислительными ресурсами (такими как серверы) с целью оптимизации производительности и минимизации времени отклика. Этот метод учитывает количество активных соединений на каждом сервере и присваивает соответствующие веса. Целью является распределение рабочей нагрузки таким образом, чтобы сбалансировать нагрузку и минимизировать время отклика.

IP Hash

Метод балансировки IP Hash использует алгоритм хэширования для определения того, какой сервер должен получить каждый из входящих пакетов. Это полезно, когда за одним IP-адресом находится несколько серверов, и вы хотите убедиться, что каждый пакет с IP-адреса клиента направляется на один и тот же сервер. Он берёт IP-адрес источника и IP-адрес назначения и создаёт уникальный хэш-ключ. Затем он используется для распределения клиента между определёнными серверами.

Это очень важный момент в случае развёртывания canary. Это позволяет нам, разработчикам, выпускать изменения для отдельной части пользователей, чтобы они могли всё протестировать и предоставить отзывы, прежде чем отправлять их в релиз.

Преимущество этого подхода в том, что он может обеспечить более высокую производительность, чем другие методы, такие как round-robin.

URL Hash

Регулировка нагрузки URL Hash также использует алгоритм хэширования для определения того, какой сервер получит каждый из запросов на основе URL.

Он также похож на метод балансировки IP Hash, но разница в том, что мы хэшируем конкретные URL, а не IP. Это гарантирует, что все запросы равномерно распределяются между серверами, обеспечивая лучшую производительность и надёжность.

HaProxy

Алгоритмы балансировки нагрузки.

- Указание серверов в разделе бэкенда позволяет HAProxy использовать эти серверы для балансировки нагрузки в соответствии с алгоритмом циклического перебора, когда это возможно.
- Алгоритмы балансировки используются для определения того, на какой сервер в бэкенде передается каждое соединение. Вот некоторые из полезных опций:
- **Roundrobin:** каждый сервер используется по очереди в соответствии со своим весом. Это самый плавный и честный алгоритм, когда время обработки серверами остается равномерно распределенным. Этот алгоритм является динамическим, что позволяет регулировать вес сервера на лету.
- Leastconn: выбирается сервер с наименьшим количеством соединений. Циклический перебор выполняется между серверами с одинаковой нагрузкой. Использование этого алгоритма рекомендуется для длинных сеансов, таких как LDAP, SQL, TSE и т. д., но он не очень подходит для коротких сеансов, таких как HTTP.
- **First**: первый сервер с доступными слотами для подключения получает соединение. Серверы выбираются от самого низкого числового идентификатора до самого высокого, который по умолчанию соответствует положению сервера в ферме. Как только сервер достигает значения maxconn, используется следующий сервер.
- **Source:** IP-адрес источника хешируется и делится на общий вес запущенных серверов, чтобы определить, какой сервер будет получать запрос. Таким образом, один и тот же IP-адрес клиента будет всегда доставаться одному и тому же серверу, в то время как серверы остаются неизменными.

Справочные материалы

- Сети docker: https://habr.com/ru/post/333874/
- Балансировка нагрузки: https://timeweb.com/ru/community/articles/algoritmy-i-metody-raspredeleniya-nagruzki-na-server
- Tier: https://habr.com/ru/company/cloud_mts/blog/332864/
- Docker compose yaml syntax: https://docs.docker.com/compose/compose-file-v3/
- Nginx:
 - https://habr.com/ru/company/first/blog/683870/
 - https://nginx.org/en/docs/http/load_balancing.html
- HaProxy: https://www.haproxy.com/blog/haproxy-configuration-basics-load-balance-your-servers/