

EVDOKIMOVA DARYA 21205
HAProxy SERVER aka ОТКАЗОУСТОЙЧИВЫЙ КЛАСТЕР

РЕКОМЕНДАЦИЯ: Всё делать с использованием мобильного интернета (короче говоря, телефон - точка доступа). Иначе сделаете из-под домашнего роутера, придёте в нгу - и там все слетит, потому что ip адреса другие, сеть другая. А с мобилой всегда одни и те же.

Здесь рассказываю про идею ha кластера. Все команды - делать по инструкции [из группы в вк](#). При возникновении проблем - гугол + может какие-то советы из моего дока помогут.

!!!ATTENTION!!!: все команды делать из-под режима `'su -l'`.

Оборудование: VirtualBox (я делала на ubuntu, вот [инструкция](#) по скачиванию);
3 сервера:

- Node1 на [debian](#) (amd64)
- Node2 на [debian](#) (amd64)
- DNSServak на [ubuntu server](#) (потому что почему-то с debian пошло что-то не так, но может у вас норм будет).

Тип - подключения - сетевой мост(!!!) (иначе у нас 3 машины будут отображаться на какой-то один адрес, который потом принимает/отдает данные).

При клонировании установленной первой машины генерация новых mac-адресов - обязательна!!! (иначе все 3 машины будут с одним и тем же ip-адресом, а нам этого не надо).

Установка машины (если чо) - как в инструкции (пункт "с. Установка OS").

Важно! Нужно прописать статические адреса на узлах (Node1 and Node2). [А зачем?](#)

ТЕРМИНОЛОГИЯ

- Кластер — группа серверов (вычислительных единиц), объединенных каналами связи.
- Отказоустойчивость (Fault Tolerance, FT) — способность системы к дальнейшей работе после выхода из строя какого-либо её элемента.
- Высокая доступность (High Availability, HA) — в случае выхода из строя узла пользователь какое-то время не будет получать услугу, однако восстановление системы произойдёт автоматически; время простоя минимизируется.
- КВД — кластер высокой доступности, HA-кластер.

Что значит "[кластер высокой доступности](#)"?

[Принцип работы ha cluster](#).

[Про corosync and pacemaker](#)

- **Corosync** — программный продукт, который позволяет создавать единый кластер из нескольких аппаратных или виртуальных серверов. Corosync отслеживает и передает состояние всех участников (нод) в кластере. Этот продукт позволяет:
 - мониторить статус приложений;
 - оповещать приложения о смене активной ноды в кластере;
 - отправлять идентичные сообщения процессам на всех нодах;
 - предоставлять доступ к общей базе данных с конфигурацией и статистикой;
 - отправлять уведомления об изменениях, произведенных в базе.
-
- **Pacemaker** - менеджер ресурсов кластера. Он позволяет использовать службы и объекты в рамках одного кластера из двух или более нод. Вот вкратце его достоинства:
 - позволяет находить и устранять сбои на уровне нод и служб;
 - не зависит от подсистемы хранения: можем забыть общий накопитель, как страшный сон;
 - не зависит от типов ресурсов: все, что можно прописать в скрипты, можно кластеризовать;
 - поддерживает STONITH (Shoot-The-Other-Node-In-The-Head), то есть умершая нода изолируется и запросы к ней не поступают, пока нода не отправит сообщение о том, что она снова в рабочем состоянии;
 - поддерживает кворумные и ресурсозависимые кластеры любого размера;
 - поддерживает практически любую избыточную конфигурацию;
 - может автоматически реплицировать конфиг на все узлы кластера — не придется править все вручную;
 - можно задать порядок запуска ресурсов, а также их совместимость на одном узле;
 - поддерживает расширенные типы ресурсов: клоны (когда ресурс запущен на множестве узлов) и дополнительные состояния (master/slave и подобное) — актуально для СУБД (MySQL, MariaDB, PostgreSQL, Oracle);
 - имеет единую кластерную оболочку CRM с поддержкой скриптов.

ИДЕЯ

(*Здесь рассказываю про идею на кластера. Все команды - делать по инструкции. При возникновении проблем - гугол + может какие-то советы из моего дока помогут. Ну и еще уточнения команд будут, зачем что делали и тд).

Идея заключается в том, что с помощью Corosync мы построим кластер (грубо говоря, “объединим” хосты), а следить за его состоянием будем с помощью Pacemaker (его главная задача — достижение максимальной доступности ресурсов, которыми он управляет, и защита их от сбоев).

[Пример](#) из жизни сисадмина, поможет понять, зачем нужен ha cluster.

ТЕПЕРЬ ДЕТАЛЬКИ

Хороший [сайт](#) с объяснениями про racemaker and corosync.

Про corosync: Смотреть презентацию в обсуждении задания к лабе.

Вот [здесь](#) хорошо написано.

Corosync uses the totem protocol for "heartbeat" like monitoring of the other node's health.

A token is passed around to each node, the node does some work (like acknowledge old messages, send new ones), and then it passes the token on to the next node. This goes around and around all the time. Should a node not pass it's token on after a short time-out period, the token is declared lost, an error count goes up and a new token is sent. If too many tokens are lost in a row, the node is declared lost/dead.

Once the node is declared lost, the remaining nodes reform a new cluster. If enough nodes are left to form quorum, then the new cluster will continue to provide services. In two-node clusters, quorum is disabled so each node can work on it's own.

В директиве totem изменяем (либо если его нет, то создаем) interface на следующий:

```
interface {  
    ringnumber: 0  
    bindnetaddr: 192.168.106.0  
    broadcast: yes  
}
```

Bindnetaddr – адрес сети нашего кластера (зависит от ip-адреса нашей виртуальной машины)

В директиве nodelist добавляем следующее:

```
node {  
    name: node1  
    nodeid: 1  
    ring0_addr: 192.168.106.104  
}  
node {  
    name: node2  
    nodeid: 2  
    ring0_addr: 192.168.106.105  
}
```

ring0_addr – это ip-адреса 1 и 2 узла

*Все настройки corosync -все по инструкции.

!При копировании ключа authkey возможна ошибка permission denied.

Совет: скопировать из-под рута (с Node1) на user (Node2) в корневую папку. А потом уже в нужную (/etc/corosync/).

Once the node is declared lost, the remaining nodes reform a new cluster. If enough nodes are left to form quorum, then the new cluster will continue to provide services. In two-node clusters, quorum is disabled so each node can work on it's own. (Это как раз наш случай)

*тут плавно перейдем к настройке Pacemaker

Читаем про [кворум](#).

Что такое кворум? Говорят, что кластер имеет кворум при достаточном количестве «живых» узлов кластера. Достаточность количества «живых» узлов определяется по формуле: $n > N/2$, где n – количество живых узлов, N – общее количество узлов в кластере.

Читаем про команды pacemaker-a.

В режиме конфигурации “crm configure” пишем следующие команды:

- **`property no-quorum-policy=ignore`**

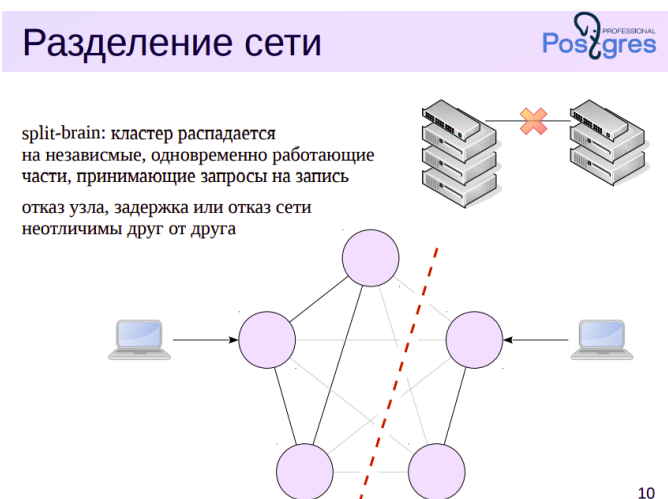
Это означает, что если если нода может достучаться более, чем до 50% нод, то нужно убить другую группу нод. В данном случае при отказе 1 ноды вторая не будет работать, т.к. не набрано 50%.

В кластере, состоящем из двух узлов, при сбое на одном из 2-х узлов не будет кворума. По умолчанию, если нет кворума, Pacemaker останавливает ресурсы. Чтобы этого избежать, нужно при настройке Pacemaker указать ему, чтобы наличие или отсутствие кворума не учитывалось.

(*Кворум - количество нод).

ВОПРОС: зачем нам вообще нужен кворум? Почему без него нельзя обойтись.

ОТВЕТ: Нужно избежать split brain - ситуация, когда кластер распадается на



независимые, одновременно работающие части, принимающие запросы на запись. Возникает рассогласование данных, и при восстановлении данные одной из групп будут вынуждено потеряны. Читаем [слайд номер 10](#). Еще [вот эта статья](#) полезная. Оттуда сейчас будет пример и объяснение.

Итак.

Важно понимать, что смерть узла системы и потеря с ним

связи не эквивалентны.

Разница проста — умерший узел не может совершать действий, которые «по неведению» окажутся деструктивными по отношению к остальной системе. Всегда можно устроить так, чтобы он умирал навсегда, или приходя в себя первым делом консультировался с остальной системой, чтобы ненароком ничего не сломать.

При разрыве связи все печальнее. Узел не может позволить себе умереть, потому что возможно только он и обеспечивает работоспособность (без связи с другими узлами невозможно понять, есть ли еще кто-то живой). Не может он и синхронизировать свою работу с остальными — связи-то нет. Остается действовать вслепую, на свой страх и риск.

Собственно, пример.

Допустим, в нашей системе ровно 2 совершенно одинаковых узла — А и В. Каждый из них хранит копию данных второго и может независимо обрабатывать запросы извне. Каждый, обрабатывая запрос, уведомляет второго об изменениях, чтобы данные оставались согласованы.

Вариант 1: узел А умирает.

Система продолжает работать как ни в чем не бывало — В продолжает обрабатывать запросы. Когда А придет в чувство, он первым делом синхронизируется с В и они вдвоем продолжат работать дальше. Ни доступность, ни согласованность не страдают.

Вариант 2: А и В живы, но связь между ними прервана.

При этом каждый из них продолжает принимать запросы извне, но не может уведомить второго об изменениях. Для каждого узла все выглядит так, будто второй узел умер и он действует в одиночку. Эту ситуацию часто называют «split-brain» — мозг разделился на два полушария, каждое из которых считает себя единственным хозяином ситуации. Система заболела шизофренией.

Если в этот момент на А был обработан запрос на удаление некой записи R, а на В был обработан запрос на модификацию той же самой записи, то данные стали несогласованы. Когда связь между А и В восстановится, при синхронизации всплывет конфликт — удалить R или оставить модифицированную версию? Тут можно выкручиваться разными стратегиями разрешения конфликтов, но consistency мы уже потеряли.

ТО ЕСТЬ решение о том, что делать с записью R определяется БОЛЬШИНСТВОМ голосов (The cluster quorum is the majority of voting nodes in the active cluster membership plus a witness vote). И именно отсюда и вытекает ответ на вопрос: “Почему отключаются ресурсы именно при более 50 процентах?”. А не 30, например.

Еще [статья полезная](#). И [еще](#). [Про CAP-теорему](#) (для общего развития, так сказать).

- `property stonith-enabled=false`

STONITH = Shoot-The-Other-Node-In-The-Head

When pacemaker is told that membership has changed because a node died, it looks to see what services might have been lost. Once it knows what was lost, it looks at the rules it's been given and decides what to do.

Generally, the first thing it does is "stonith" (aka "fence") the lost node. This is a process where the lost node is powered off, called power fencing, or cut off from the network/storage, called fabric fencing. In either case, the idea is to make sure that the lost node is in a known state. If this is skipped, the node could recover later and try to provide cluster services, not having realized that it was removed from the cluster. This could cause problems from confusing switches to corrupting data.

В нашем случае каждая нода будет пытаться сделать STONITH другой ноде. Ничего работать не будет, потому что все ноды умрут.

- primitive VIP ocf:heartbeat:IPaddr2 params ip=192.168.0.11 cidr_netmask=24 op monitor interval=1s

Создаём виртуальный ip-адрес(VIP, то есть создаем ресурс); с помощью op monitor interval=1s кластер будет следить за состоянием ресурса каждую секунду. По-умолчанию кластер не обеспечивает работоспособность ресурсов, т.е. не следит после запуска, жив ли ресурс. Чтобы это происходило, нужно добавлять операцию monitor, тогда кластер будет следить за состоянием ресурса. Параметр interval этой операции - интервал с каким делать проверку.

- primitive HAP lsb:haproxy op monitor interval=1s

Создаём ресурс haproxy(HAP)

- colocation CLC inf: VIP HAP

Это принуждение к размещению ресурсов на одном узле или наоборот на разных узлах. При отказе одной ноды, вторая должна иметь необходимые ресурсы для продолжения работы, поэтому необходимо разместить ресурсы на *каждой* ноде.

Иногда нужно, чтобы некоторый набор ресурсов выполнялся на одном узле, но не все эти ресурсы являются примитивами, поэтому их нельзя сгруппировать. В этом случае используется colocation - это принуждение к размещению ресурсов на одном узле или наоборот на разных узлах.

- order ORD inf: VIP HAP

Определяем порядок запуска ресурсов (сначала VIP, а потом HAP, зависящий от VIP)

- commit

Сохраняем настройки.

- /sbin/reboot

Перезапускаем машину.

*Когда установим pacemaker на Node2 и напишем 'crm configure status', то будут те же данные, что и на Node1 (на которой мы изначально задавали все настройки).

ПРО ВЗАИМОСВЯЗЬ МЕЖДУ COROSYNC-ОМ И PACEMAKER-ОМ

<https://github.com/voxpupuli/puppet-corosync/issues/32>

HAPROXY

Теперь перейдем к haproxy. HAProxy - это балансировщик нагрузки.

Отличная [статья](#). [Документация](#).

Как работает [балансировка](#)?

Команда `'net.ipv4.ip_nonlocal_bind=1'` [зачем](#) нужна? This command allows processes to bind() to non-local IP addresses, which can be quite useful for application such as load balancer.

То есть у нас есть floating point - он же frontend - адрес 192.68.106.111 - виртуальный адрес, который мы привязываем в наш кластер.

APACHE2

Читаем вот [это](#).

- `/etc/apache2/sites-available/mysite.com.conf`

Этот каталог содержит файлы настроек для **виртуальных сетевых узлов** (Virtual Hosts) Apache2. Виртуальные сетевые узлы позволяют настраивать Apache2 на множество сайтов с отдельными конфигурациями.

- **ports.conf**: содержит инструкции, которые определяют какие TCP порты прослушивает Apache2.
- Инструкция **Listen** определяет порт, и в общем случае IP адрес, на которых Apache2 должен ожидать соединения. Если IP адрес не определен, Apache2 будет прослушивать все IP адреса, которые назначены компьютеру, где он запущен. Значение по умолчанию для **Listen** 80. Замените его на 127.0.0.1:80 чтобы Apache2 прослушивал только интерфейс внутренней петли, что сделает его недоступным из интернета; на 81 (например) для изменения порта доступа или оставьте как есть для стандартного функционирования. Эта инструкция может быть найдена и изменена в единственном файле `/etc/apache2/ports.conf`.
- 4. Инструкция **ServerName** необязательная и определяет на какой адрес FQDN ваш сайт должен отвечать. Изначальный виртуальный хост не имеет ServerName, поэтому отвечает на все запросы не соответствующие директивам ServerName других виртуальных хостов. Если вы приобрели доменное имя `ubunturocks.com` и хотите прописать его на вашем Ubuntu сервере, значение ServerName для файла настроек вашего виртуального хоста должно быть `ubunturocks.com`. Добавьте эту инструкцию в файл нового виртуального хоста, который вы создавали ранее (`/etc/apache2/sites-available/mynewsite`).
- Инструкция **DocumentRoot** определяет где Apache2 будет искать файлы, которые являются содержимым сайта. По умолчанию используется значение `/var/www`, как определено в `/etc/apache2/sites-available/default`. Если желаете, можете изменить это значение в файле сайта вашего виртуального хоста и не забудьте создать этот каталог, если необходимо!
-

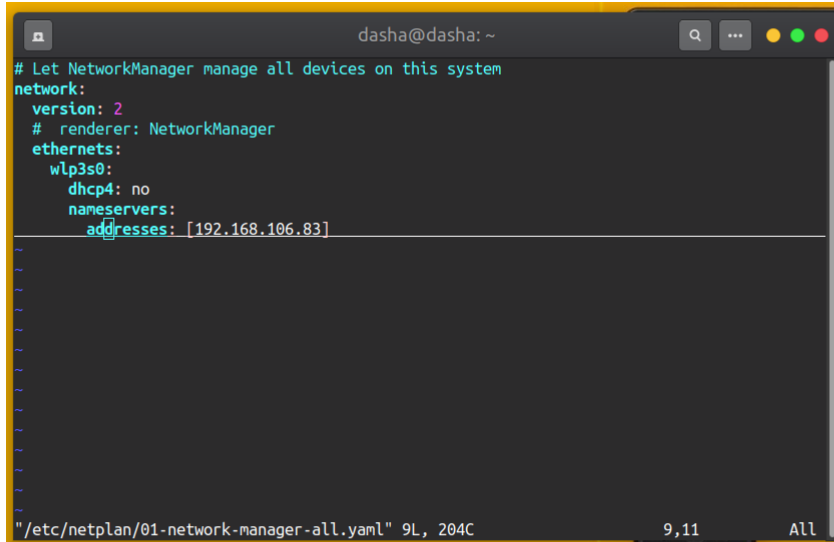
ПРО DNS

Статья на [Хабре](#).

СКРИНЫ ФАЙЛОВ

!!!ВСТАВИТЬ ИХ!!!

Вот ток для настройки того, чтоб с ноута был доступ к сайтам. *dhcp4 там может быть и yes. то работает то не работает(это заметка)

A screenshot of a terminal window with a dark background. The window title is 'dasha@dasha: ~'. The terminal displays a YAML configuration for NetworkManager. The configuration includes a comment about NetworkManager managing all devices, followed by 'network:' and 'version: 2'. Under 'ethernets:', there is a section for 'wlp3s0:' with 'dhcp4: no', 'nameservers:', and 'addresses: [192.168.106.83]'. The terminal shows a vertical scrollbar on the right side. At the bottom, the file path '/etc/netplan/01-network-manager-all.yaml' is shown along with its size '9L, 204C' and a status '9,11 All'.

```
# Let NetworkManager manage all devices on this system
network:
  version: 2
  # renderer: NetworkManager
  ethernets:
    wlp3s0:
      dhcp4: no
      nameservers:
        addresses: [192.168.106.83]
```

"/etc/netplan/01-network-manager-all.yaml" 9L, 204C 9,11 All