



Что под капотом у облачного PostgreSQL в Ozon

Дмитрий Васильев

Эксперт разработки информационных систем группы PostgreSQL DBA, Ozon

Ozon Tech 2022

PostgreSQL в OZON

- Почти каждый новый сервис имеет базу PostgreSQL а иногда и не одну (шардирование).
- Каждый месяц в OZON появляется несколько сотен PostgreSQL кластеров.
- Производить инсталляцию вручную:
 - Медленно - мешает развитию
 - Человеческий фактор
 - Рутина

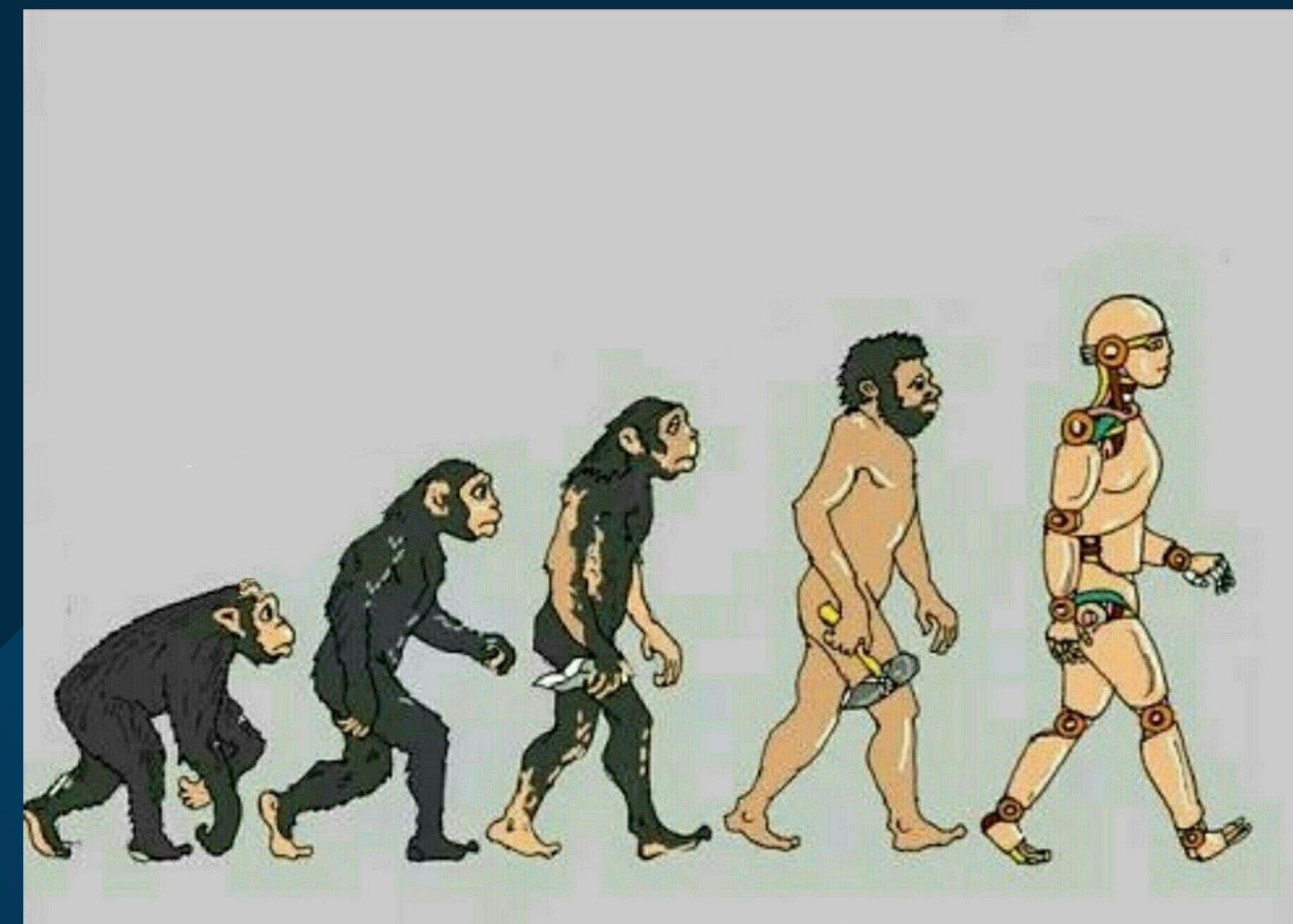




хенджоб

**МОЩНЫЙ
ЭЙПИАЙ**

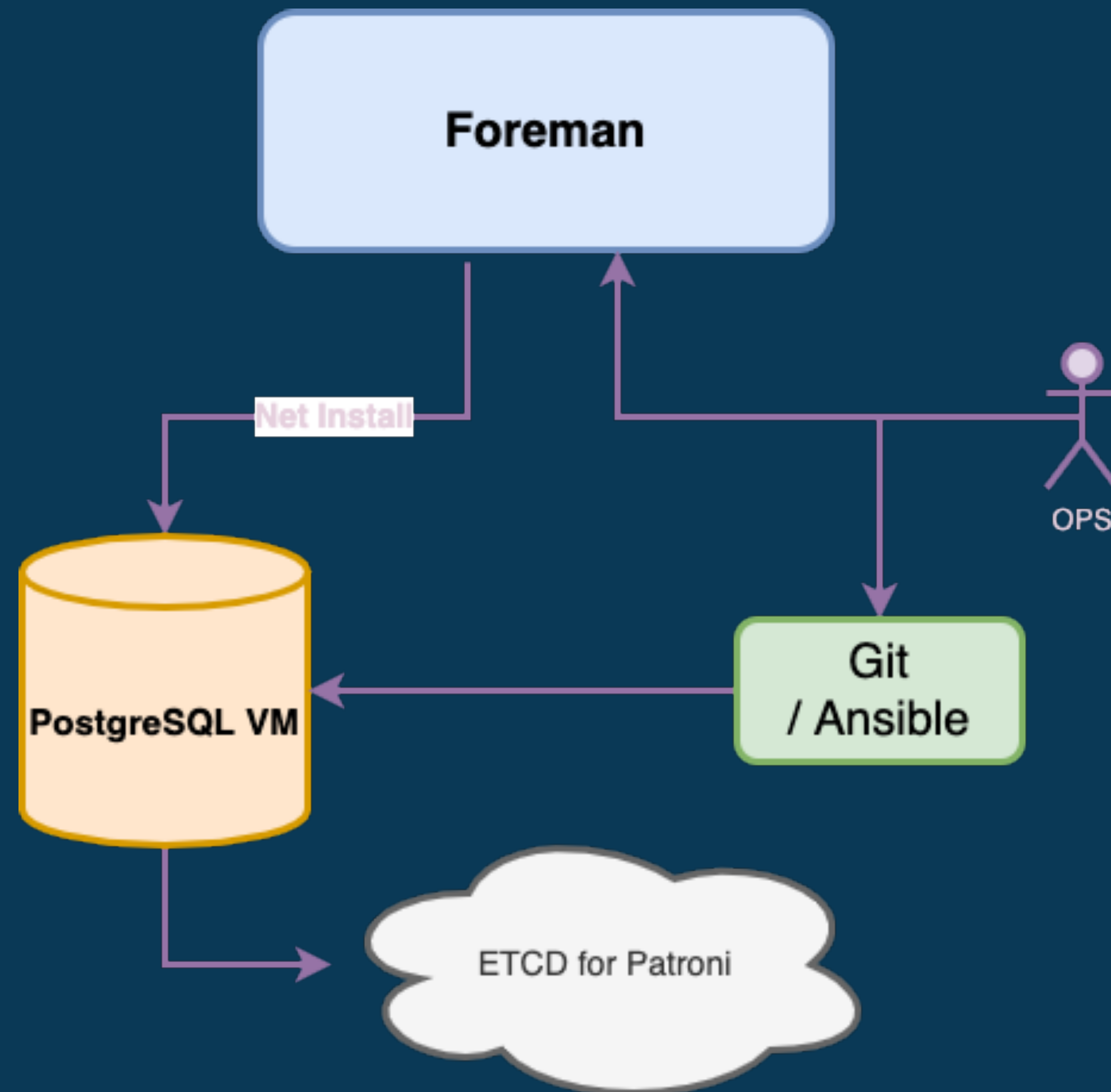
Эволюция разворачивания PostgreSQL кластера



Инсталляция PostgreSQL кластера

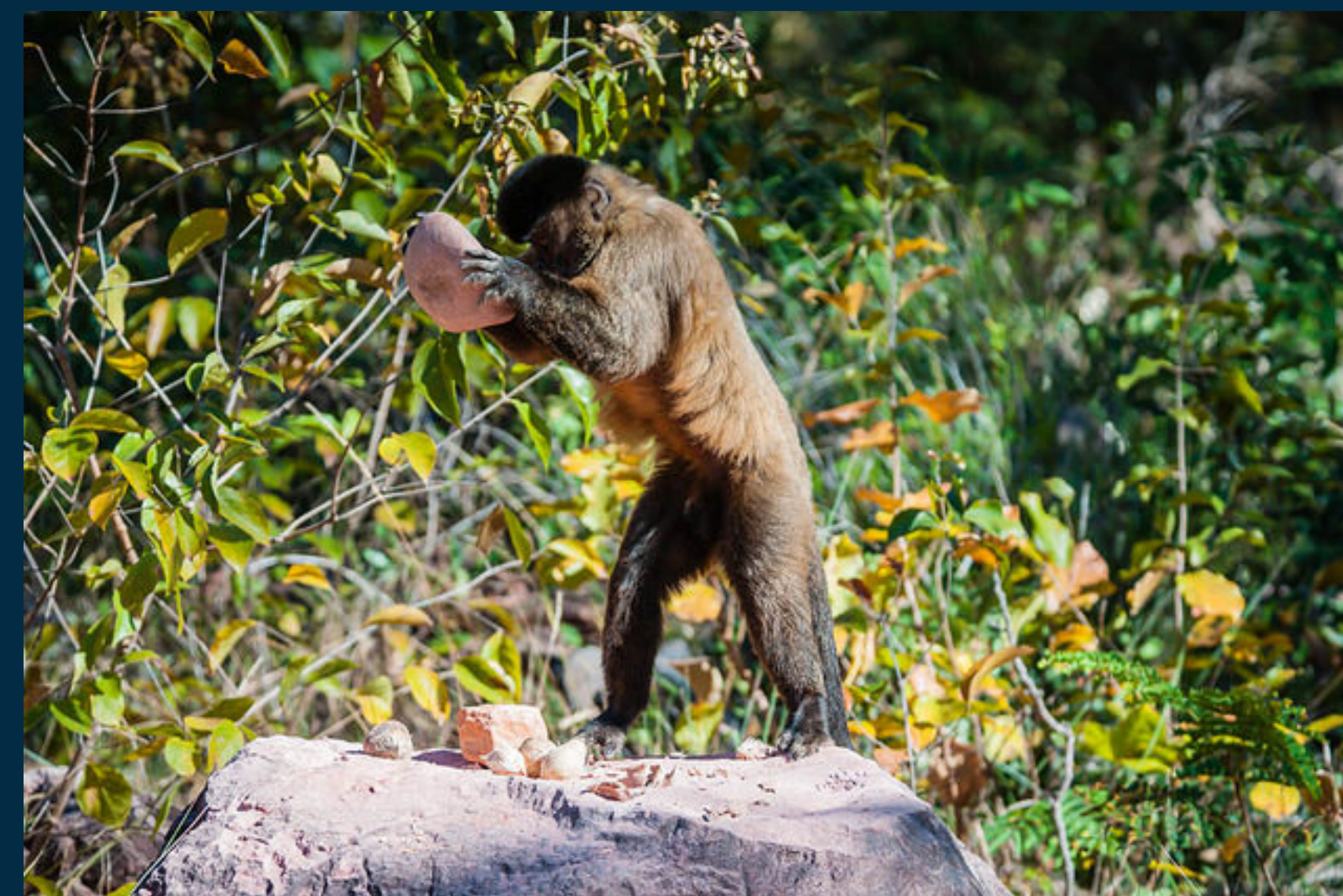
Как всё начиналось

- Первоначальная наливка VM при помощи Foreman
- GroupVars для кластера PostgreSQL через Ansible



Тупо автоматизируем ручную работу.

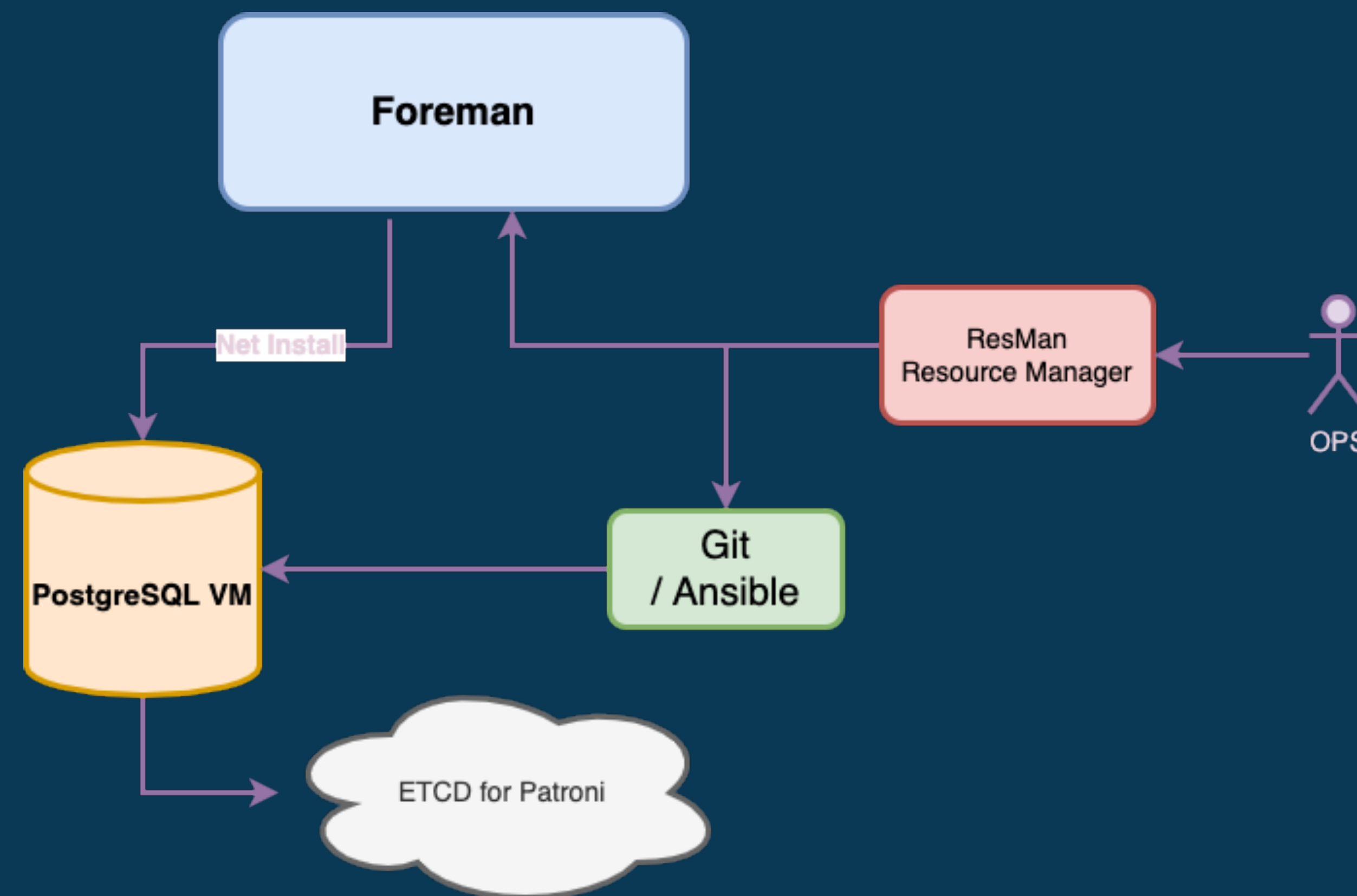
Шаг номер 1: Resource Manager.



Инсталляция PostgreSQL кластера

Первые шаги к автоматизации

- Добавлен Resource Manager
 - контролирует правильное размещение виртуалок с учетом стойки/ДЦ/выделенных ресурсов на гипервизоре
 - коммитит в Git необходимые GroupVars



Появилось несколько дата центров.

Разнородное железо под гипервизоры.

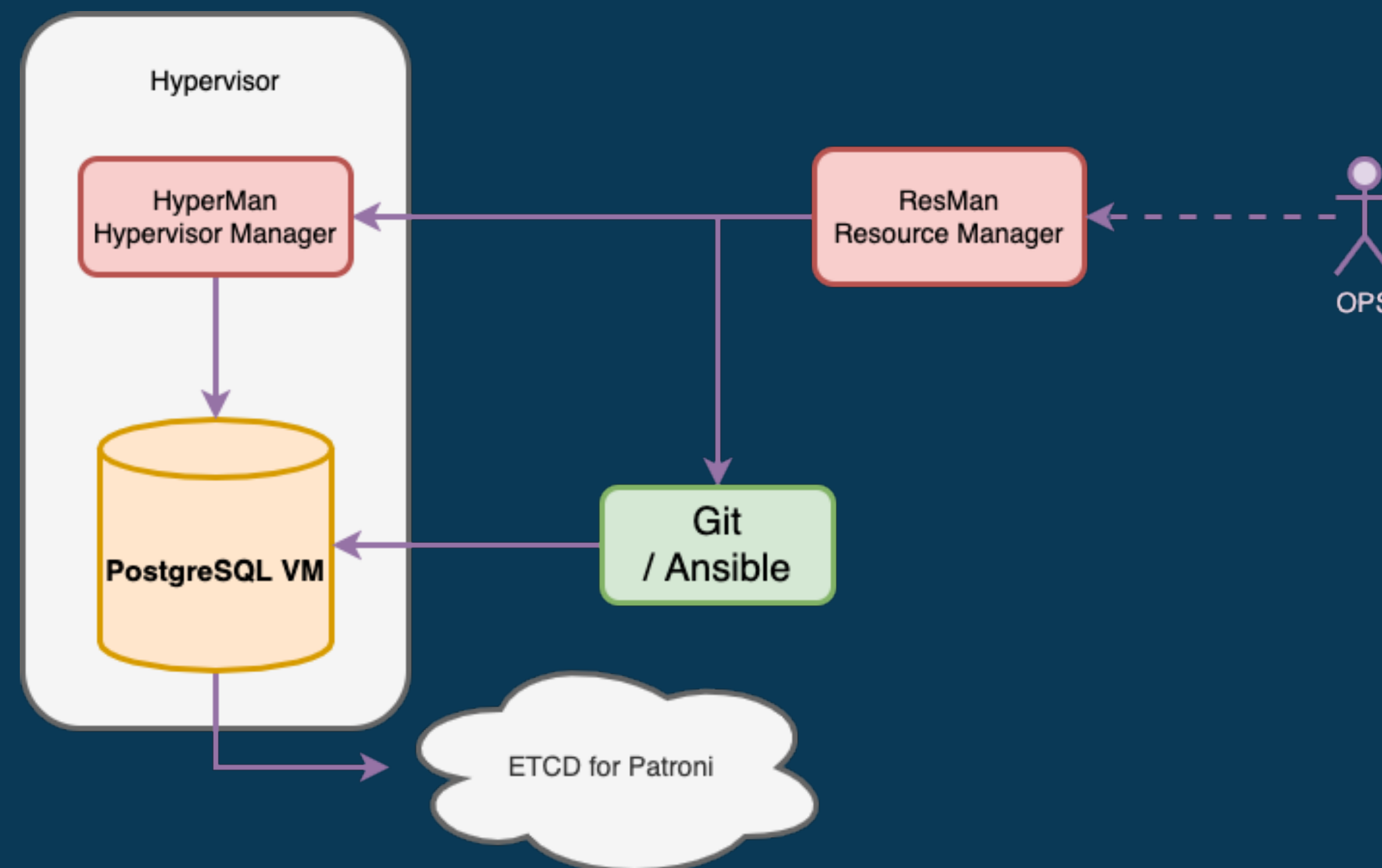
Шаг номер 2: Hyper-Man.



Инсталляция PostgreSQL кластера

Избавляемся от Foreman

- Добавлен децентрализованный Hypervisor Manager
 - настройка над Libvirt
 - использование “Gold Image”
 - онлайн ресайзы виртуальных машин через CLI
 - Интеграция с qemu-agent



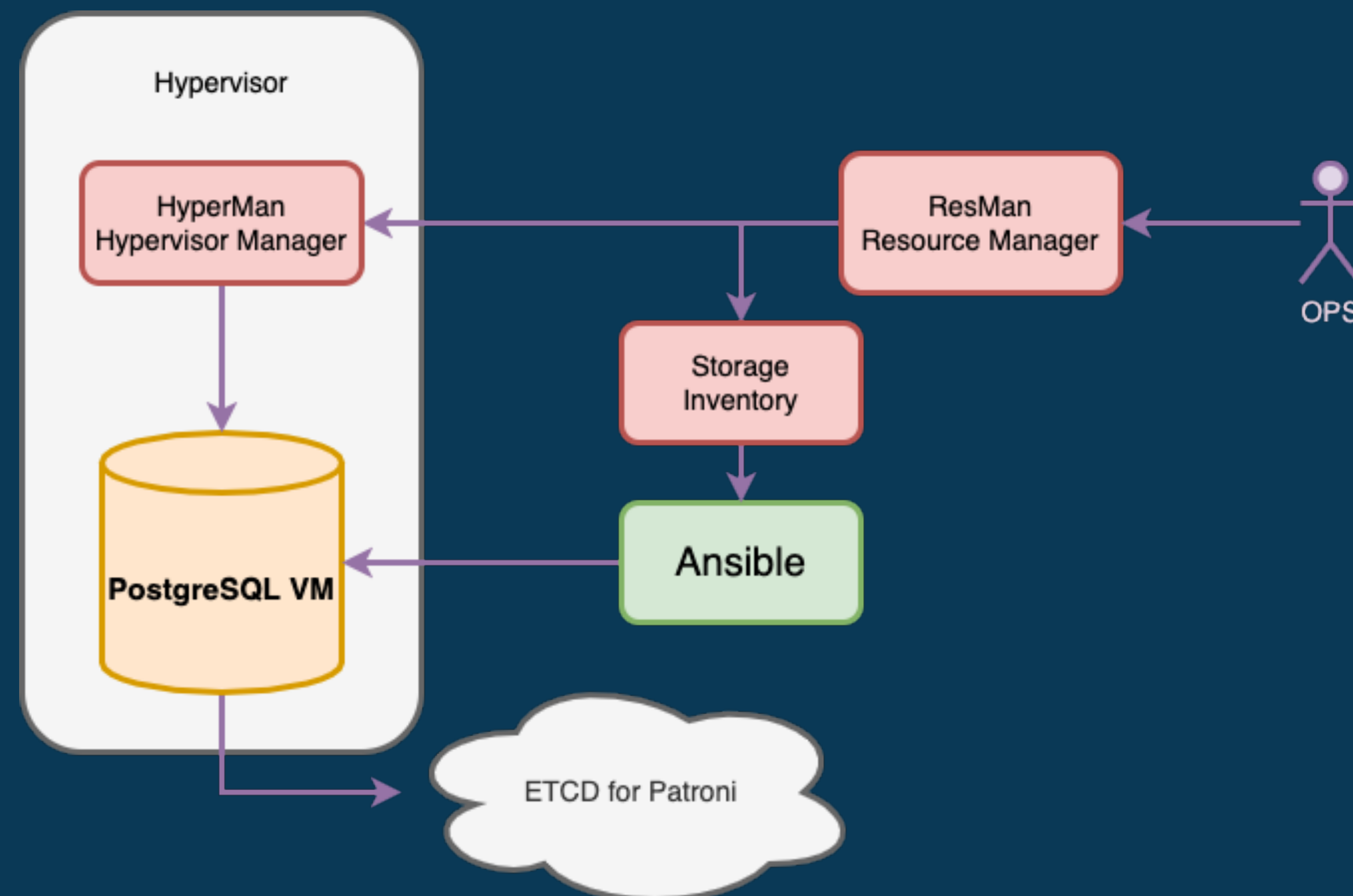
Merge approve стал бутылочным горлышком.

Шаг номер 3: Storage-Inventory.

Инсталляция PostgreSQL кластера

Избавляемся от коммитов в Git

- Добавлен Storage Inventory
 - API для хранения GroupVars
 - Аудит/История/Откаты



Боль и унижение

1. Работа по сети — очень медленно. Полный деплой всех кластеров - 1 день.
2. Ускоряемся, когда запускаем локально (`ssh <hostname> "ansible-playbook"`). Полный деплой - 2 часа.
3. Запуск очень дорогой по ресурсам (`repo cache` → `disk`, `python inventory` → `CPU`).
4. Не удовлетворяет банальным требованиям:
 1. Поставить пакет: лишний раз нужно дергать `repo cache`.
 2. Потерянные `notify` в случае падения (падение при недоступности `vault`, например).

Puppet, Ansible ... Chit!

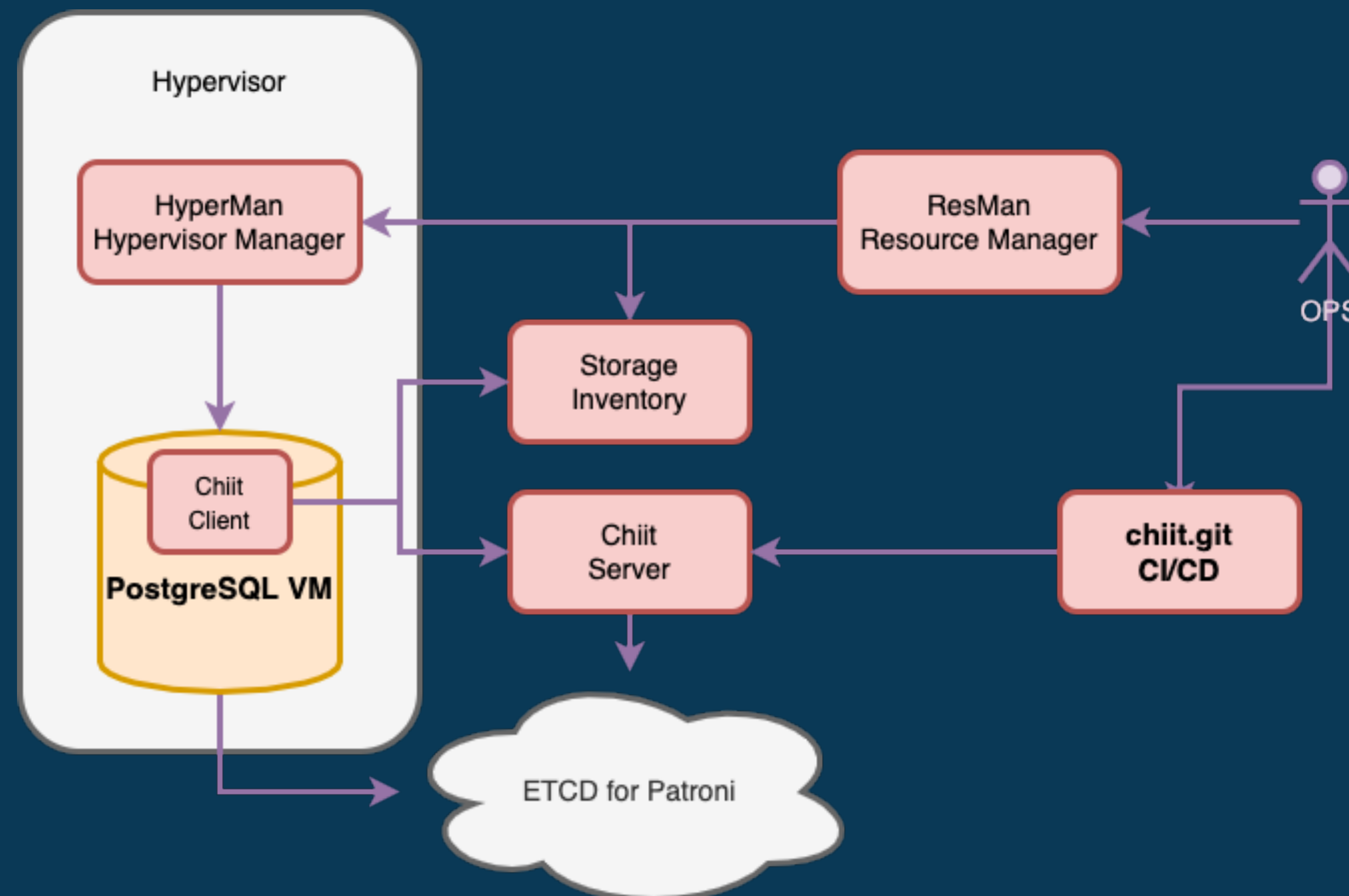
Шаг номер 4: Пишем свой SCM



Инсталляция PostgreSQL кластера

Chiit: деплой тысяч машин за 5 секунд

- Добавлен новый SCM — Chiit (Change It)
 - Клиент-серверная архитектура
 - Хранение секретов (кэш перед vault)
 - Быстрота работы
 - Клиент сохраняет отчёт о статусе и измененных ресурсах
- Интеграция с CI/CD:
 - собираем бинарь
 - тестируем через Inspec
 - загружаем бинарь в несколько Serp
 - делаем запись в Server о релизной версии



Chiit example: stop using YAML

```
systemd.NewService(ctx, &systemd.ServiceConfig{
  Name: "service-name",
  After: "after",
  ExecStart: "package-bin",
  User: "user",
  SyslogIdentifier: "service",
  MaxMemoryMb: 1024,
  CPUQuotaPercent: 50,
  Restart: systemd.ServiceRestartAlways,
  CPUAffinity: 1,
  MemoryAccounting: true,
  CPUAccounting: true,
})
apt.Install(ctx, &packages.Package{Name: "package-name",
  Version: inventory.Get(key: "package.version").String()}).Notify(
  func() {
    defers.Runner.AddRestartSystemd(ctx, name: "service-name")
  })
})
```

Инсталляция PostgreSQL кластера

Запросить Postgres DB

Имя ресурса: pg-shard-manager-1 | Среда: production

Ресурсный пул: default | Версия: По умолчанию | Количество реплик (>=3): 3

План для ресурса

- s1.micro** | Реплики: 3
2 ядра, 2 ГБ памяти, 30 ГБ диска
- s1.small-30 | Реплики: 3
4 ядра, 4 ГБ памяти, 30 ГБ диска
- s1.small-100 | Реплики: 3
4 ядра, 4 ГБ памяти, 100 ГБ диска
- s1.small-200 | Реплики: 3
4 ядра, 4 ГБ памяти, 200 ГБ диска

Квоты команды postgres_dba_team

Ресурс	Текущее значение	Лимит	Изменение
Ядра	34	36	+6
Память, ГБ	34	36	+6
Диск, ГБ	480	500	+90

Для данного плана не хватает квот по CPU, RAM и HDD

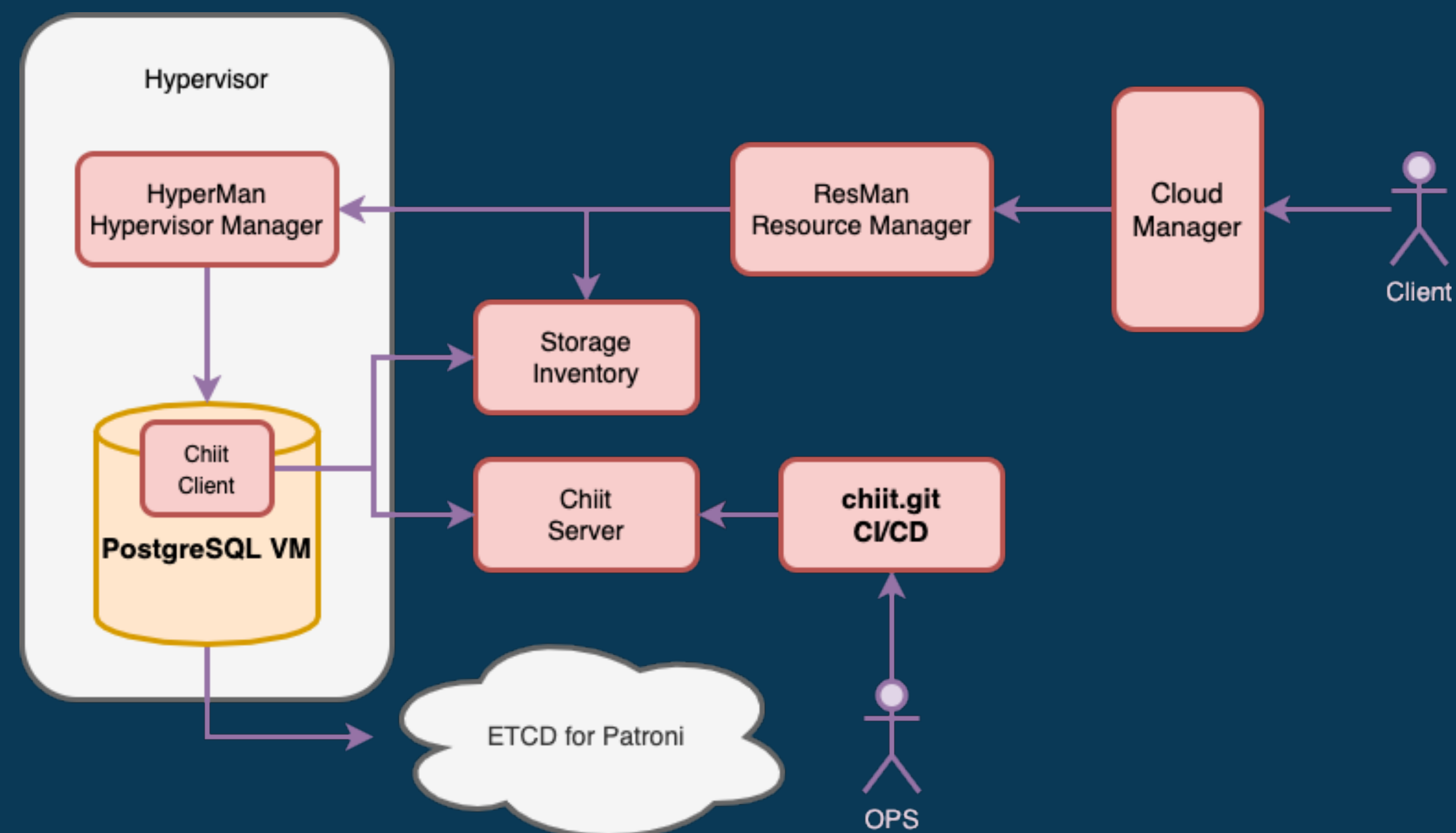
[Запросить у Platform Analytics](#)



Инсталляция PostgreSQL кластера

Поворачиваемся лицом к пользователю

- Добавлена интеграция с Cloud Manager Console
 - избавили от необходимости OPS проводить операции даже через cli
 - добавили понятие размера ресурсного пула для команды



Инсталляция PostgreSQL кластера

Редактировать ресурс Postgres DB

pg-shard-manager

Production Ресурсный пул: default Версия: 13

Внимание!
Значение диска в заказе доступно только на увеличение и не должно превышать 20% от текущего используемого.
Скейл реплик доступен только на увеличение.

Параметры ресурса

Количество реплик (≥3) 3

<input type="radio"/> s1.micro Реплики: 3	<input type="radio"/> s1.small-30 Реплики: 3
2 ядра, 2 ГБ памяти, 30 ГБ диска	4 ядра, 4 ГБ памяти, 30 ГБ диска
<input type="radio"/> s1.small-100 Реплики: 3	<input type="radio"/> s1.small-200 Реплики: 3
4 ядра, 4 ГБ памяти, 100 ГБ диска	4 ядра, 4 ГБ памяти, 200 ГБ диска
<input checked="" type="radio"/> custom Реплики: 3	
2 ядра, 2 ГБ памяти, 30 ГБ диска	

Результат:

Квоты команды postgres_dba_team

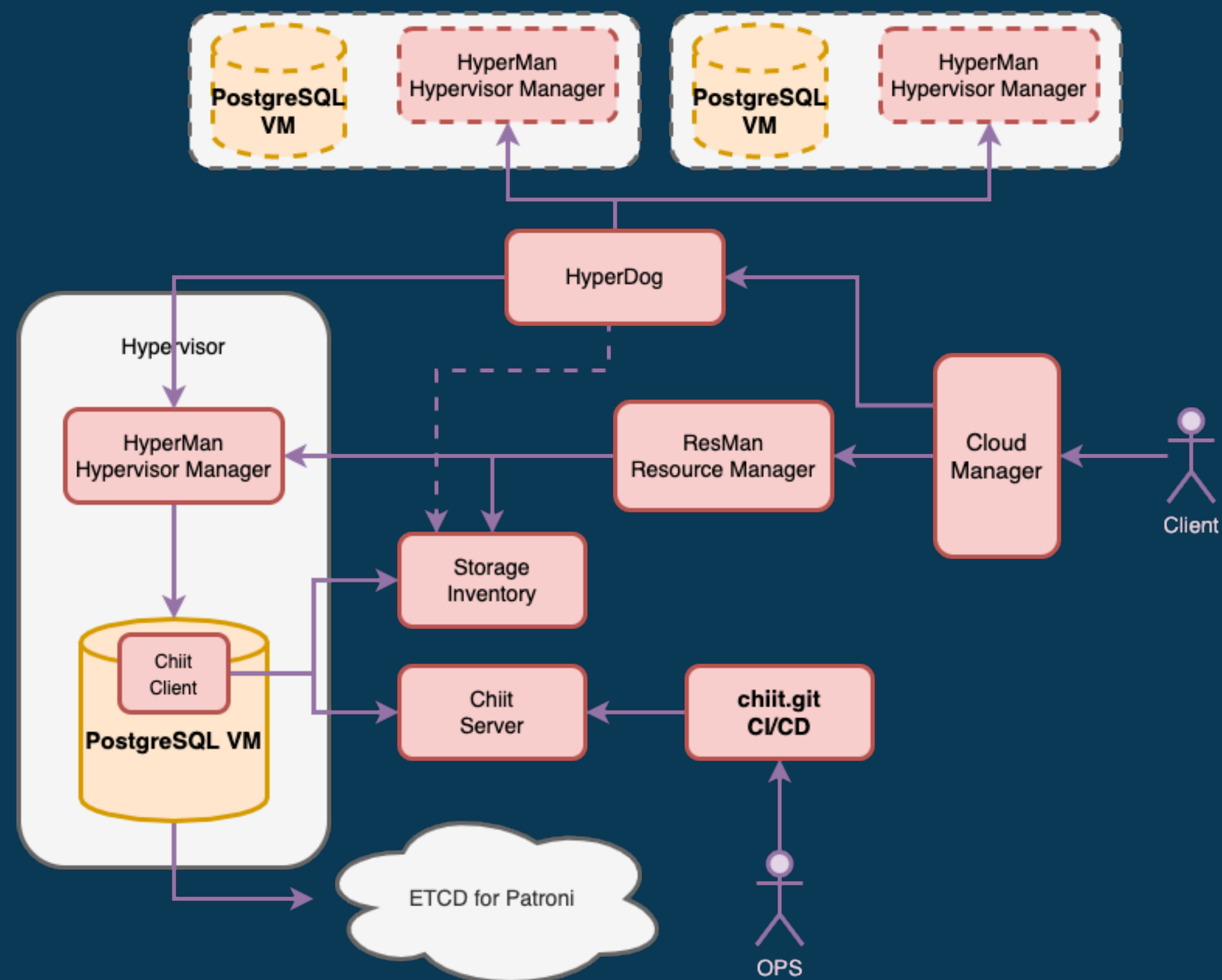
Ядра	Память, ГБ	Диск, ГБ
34 / 36	34 / 36	480 / 500
6	6	90



Инсталляция PostgreSQL кластера

Проблема роста размера кластера

- Добавлен Hyper-Dog (управляет Hyper-Man как кластером PostgreSQL).
- следит за дисковым авто-ресайзом
- производит ресайзы (CPU/Memory) кластеров



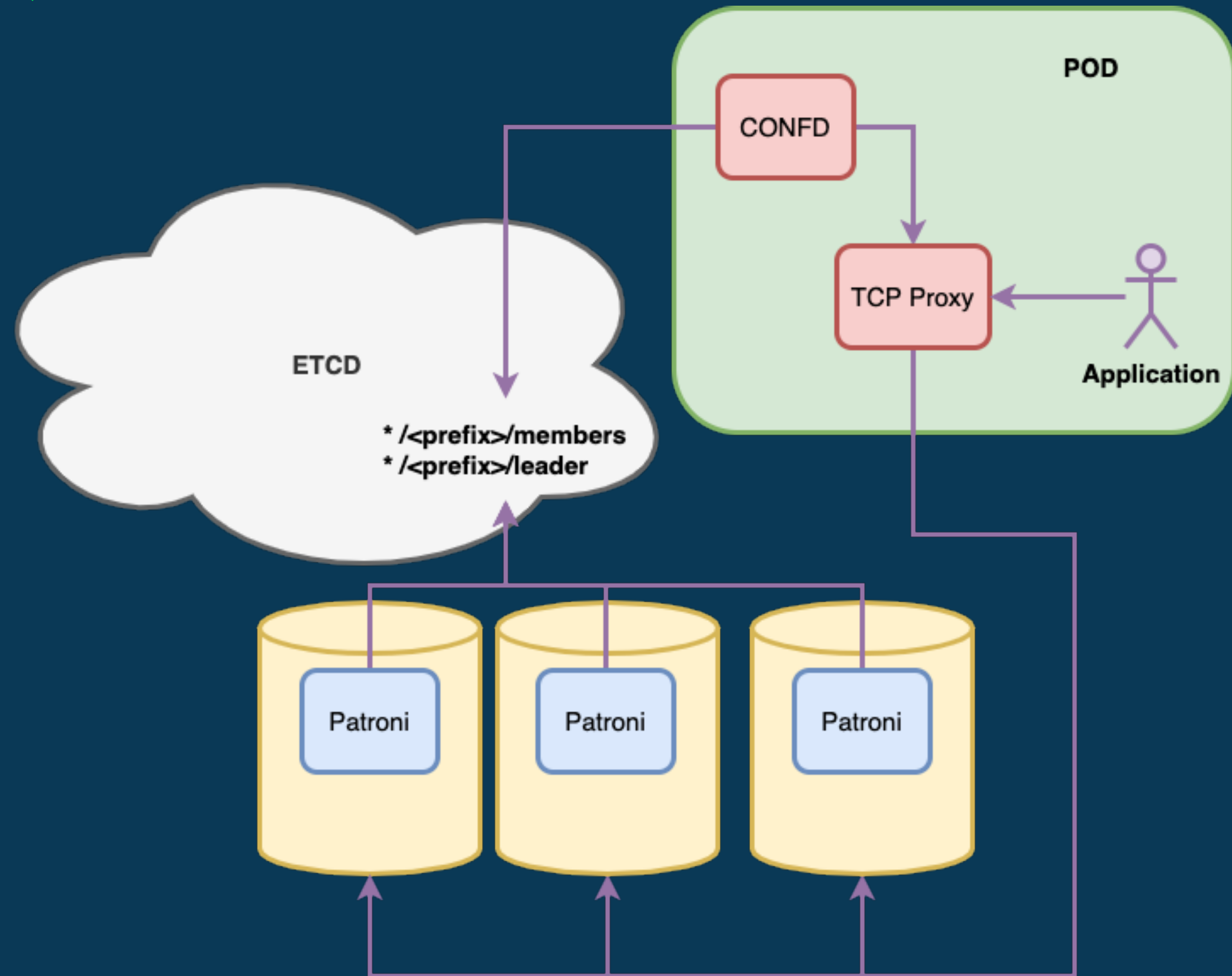
Отказоустойчивость и драйвера



Отказоустойчивость и драйвера

Подсаживаем в pod sidecar с tcp-проxy (haproxy/nginx)

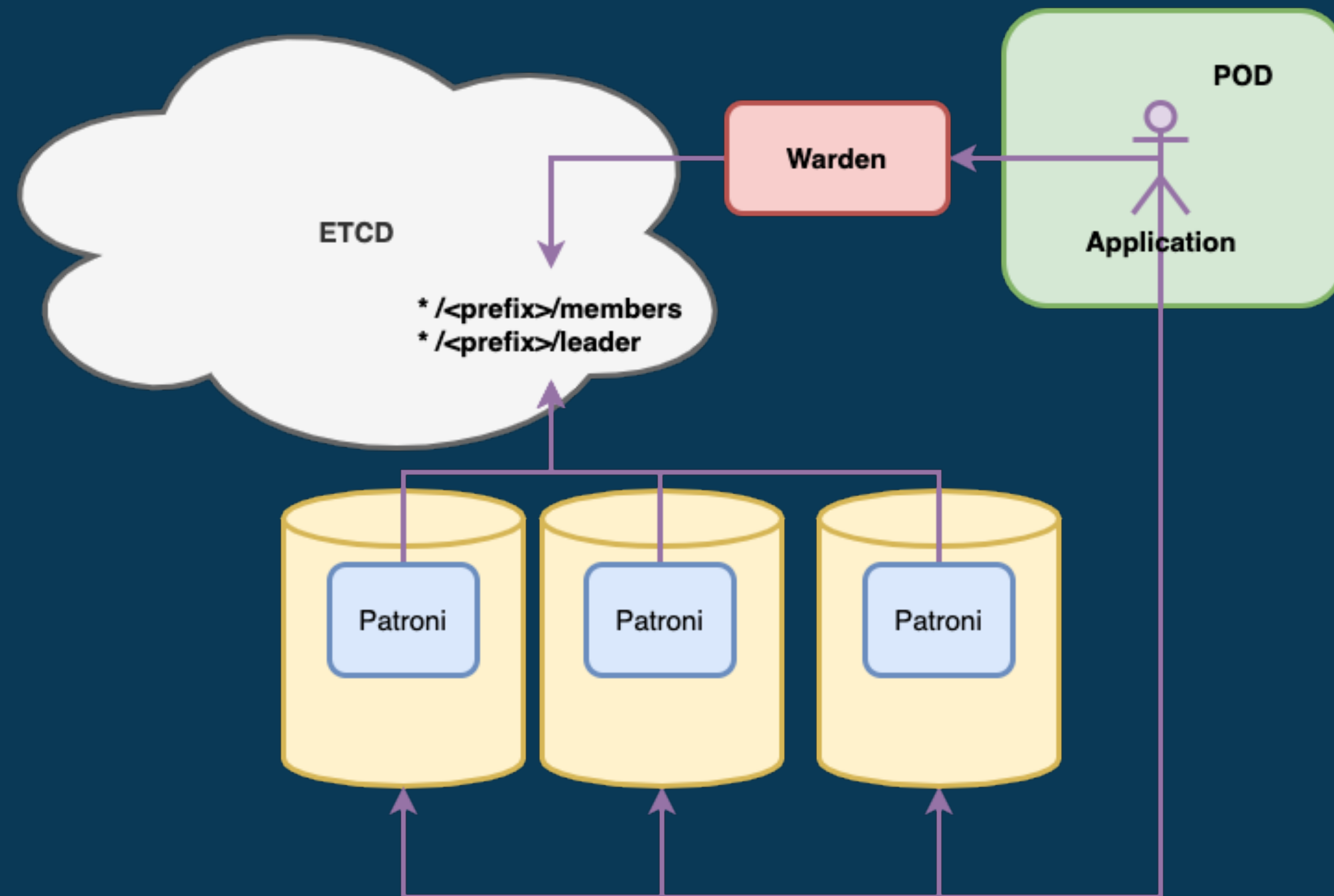
- ConfD следит за состоянием кластера в ETCD
- ConfD генерирует конфиг проxy с конечными endpoints до PostgreSQL



Отказоустойчивость и драйвера

Большое количество коннектов к ETCD, «зависание» некоторых POD'ов

- Добавили Service Discovery для PostgreSQL кластеров
 - У каждого из POD'ов единая картина мира
 - Простота разработки:
 - `Cluster.New("cluster-name").GetRole("master").Query()`
 - `Cluster.New("cluster-name").GetRole("read").Query()`
 - `Cluster.New("cluster-name").GetRole("sync").Query()`



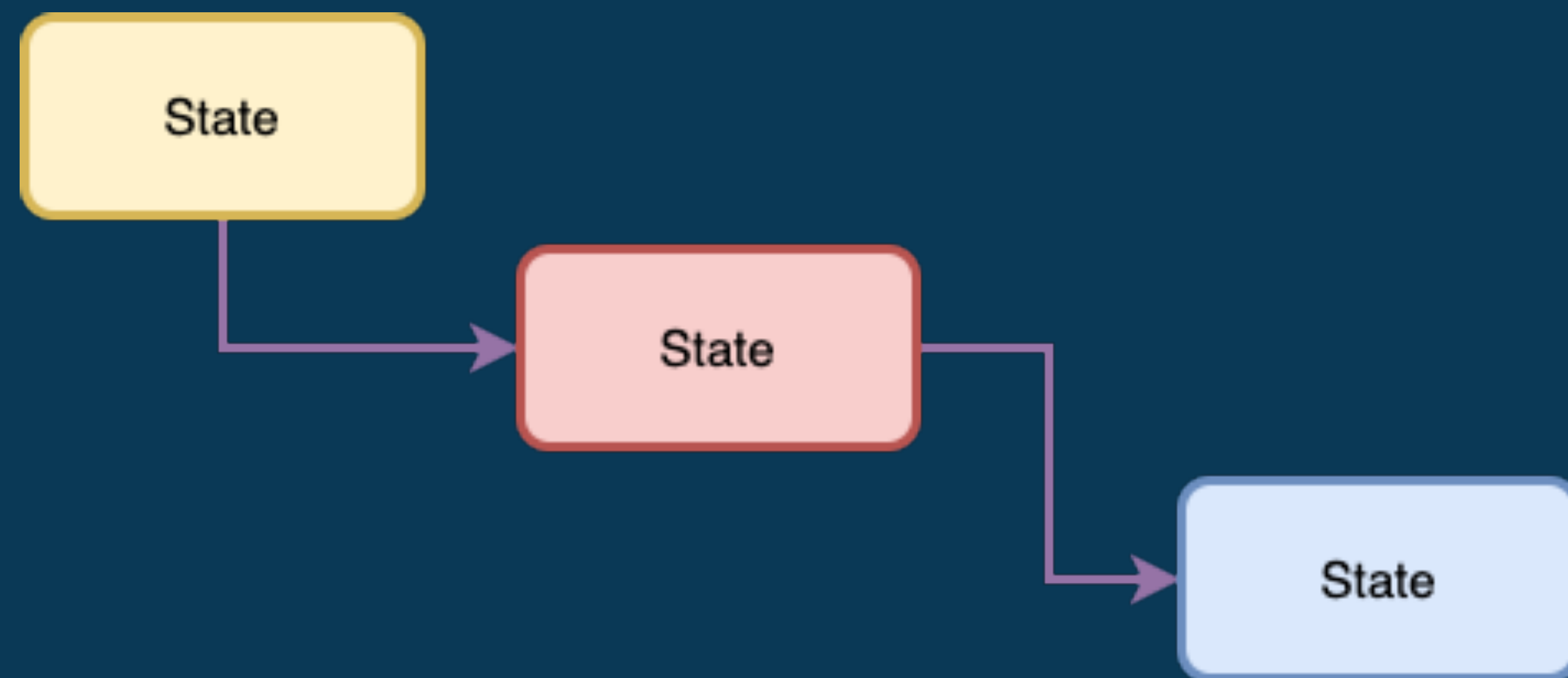
Изменение размера кластера “на лету”



Изменение кластера «на лету»

“Safe scripts”

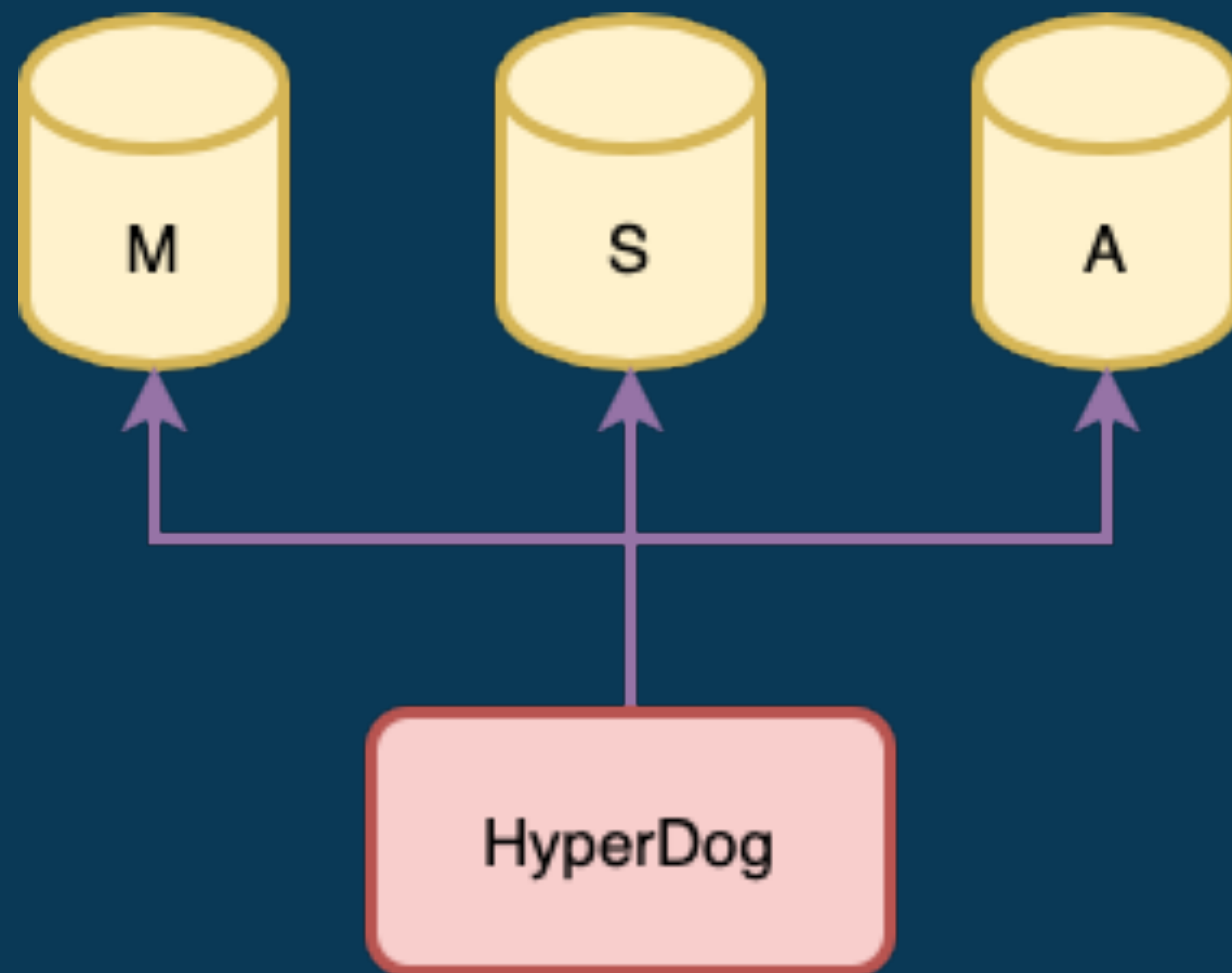
- Проект, который безопасно позволяет перезагрузить машину из кластера: `master/sync/async`
- Проверяет на целостность кластер (3 ноды)
- Делает вывод из кластера при помощи `tag NoLoadBalance` через `Service Discovery`
- Остановку PostgreSQL учитывая логическую репликацию



Изменение кластера «на лету»

Меняем ram/количество CPU

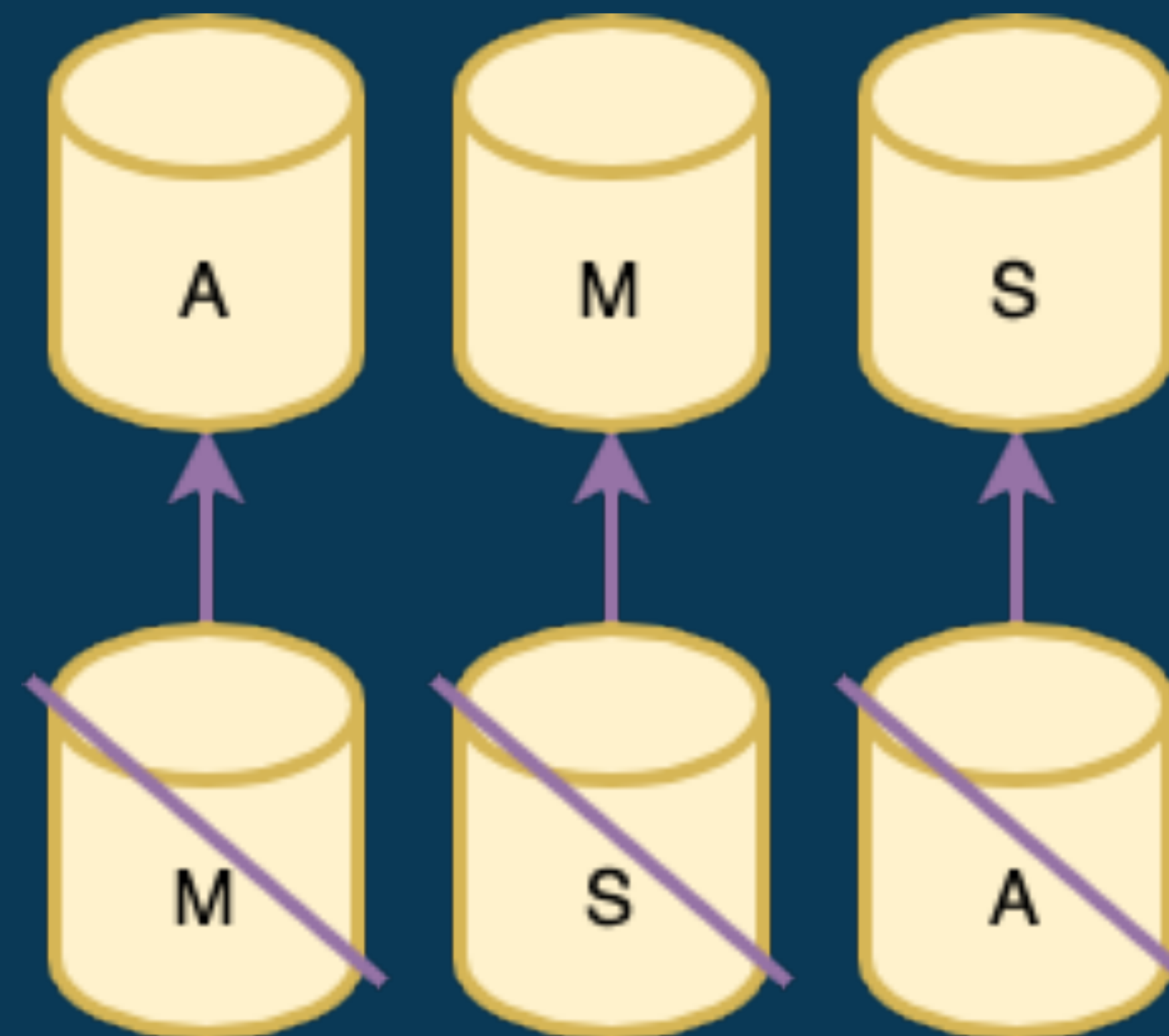
- Делаем изменения последовательно:
 - рестарт Async-реплик с новыми размерами
 - делаем switchover sync-реплики: ждем появления sync-реплики, рестарт с новыми размерами
 - делаем switchover master, рестарт с новыми параметрами



Изменение кластера «на лету»

Сжатие дисков, *shrink*.

- Определяем фактический размер базы
- Пересоздаем со *switchover async, sync, master*



Функциональность облака

Что в итоге мы умеем?

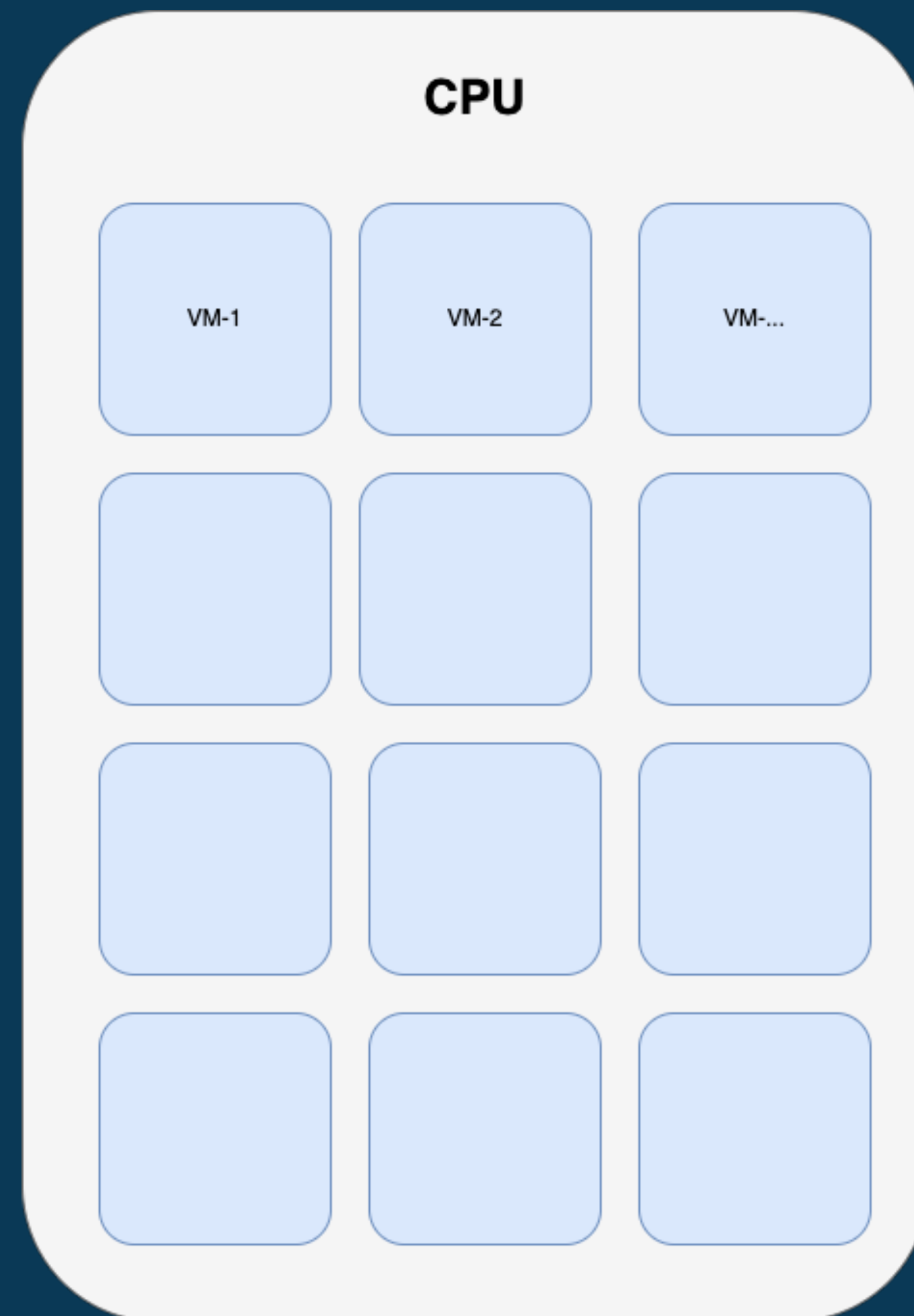
- CPU, RAM x2 без switchover.
- CPU, RAM - без аффекта для приложения.
- Disk вверх и вниз.

Тюнинг виртуальных машин

Тюнинг CPU

Первоначально общий пул CPU

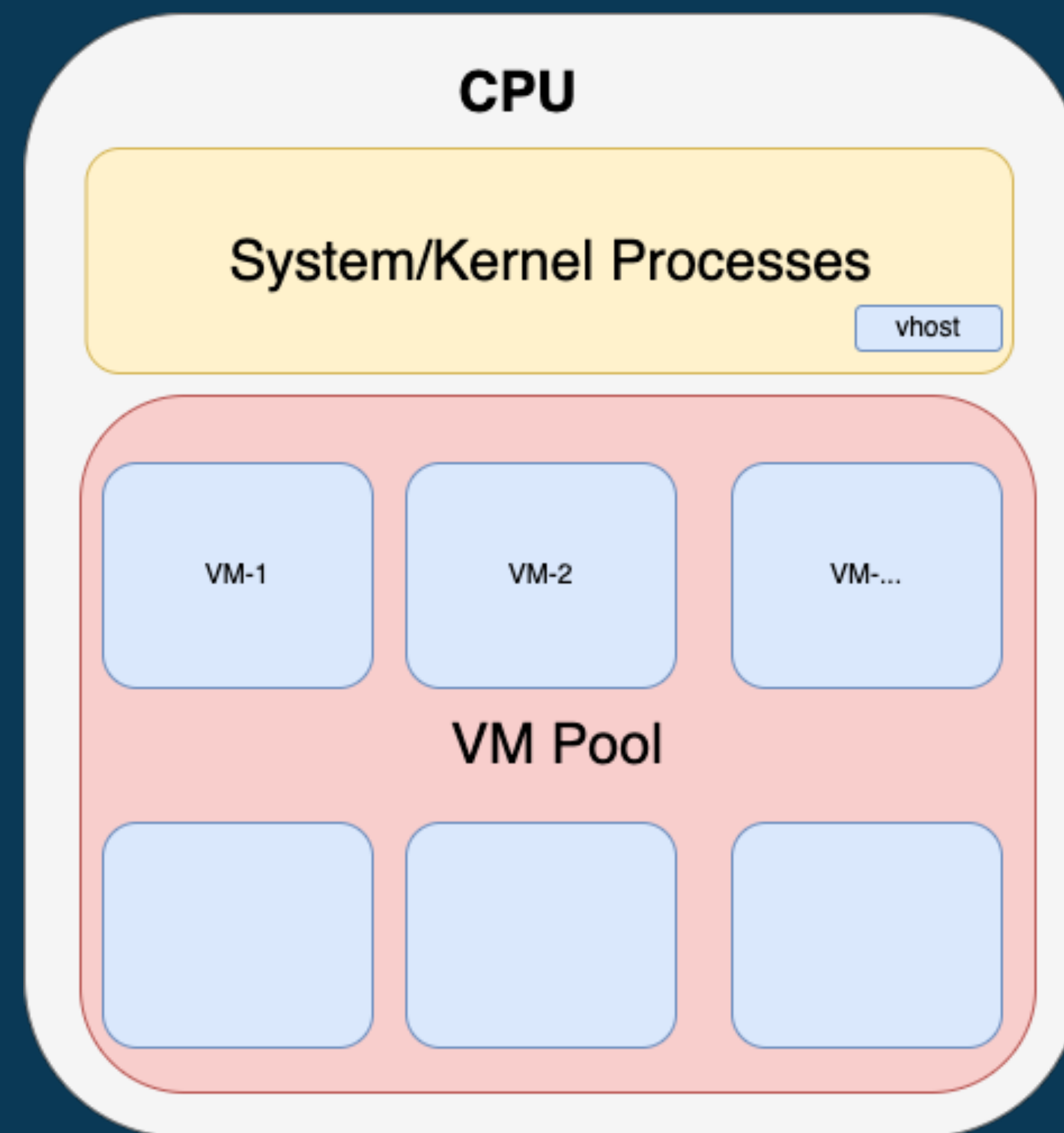
- Все системные процессы и виртуальные машины делятся всеми CPU
- Есть оверселлинг, виртуальные машины мешают друг-другу



Тюнинг CPU

Системный пул и пул для виртуальных машин: лечим сеть

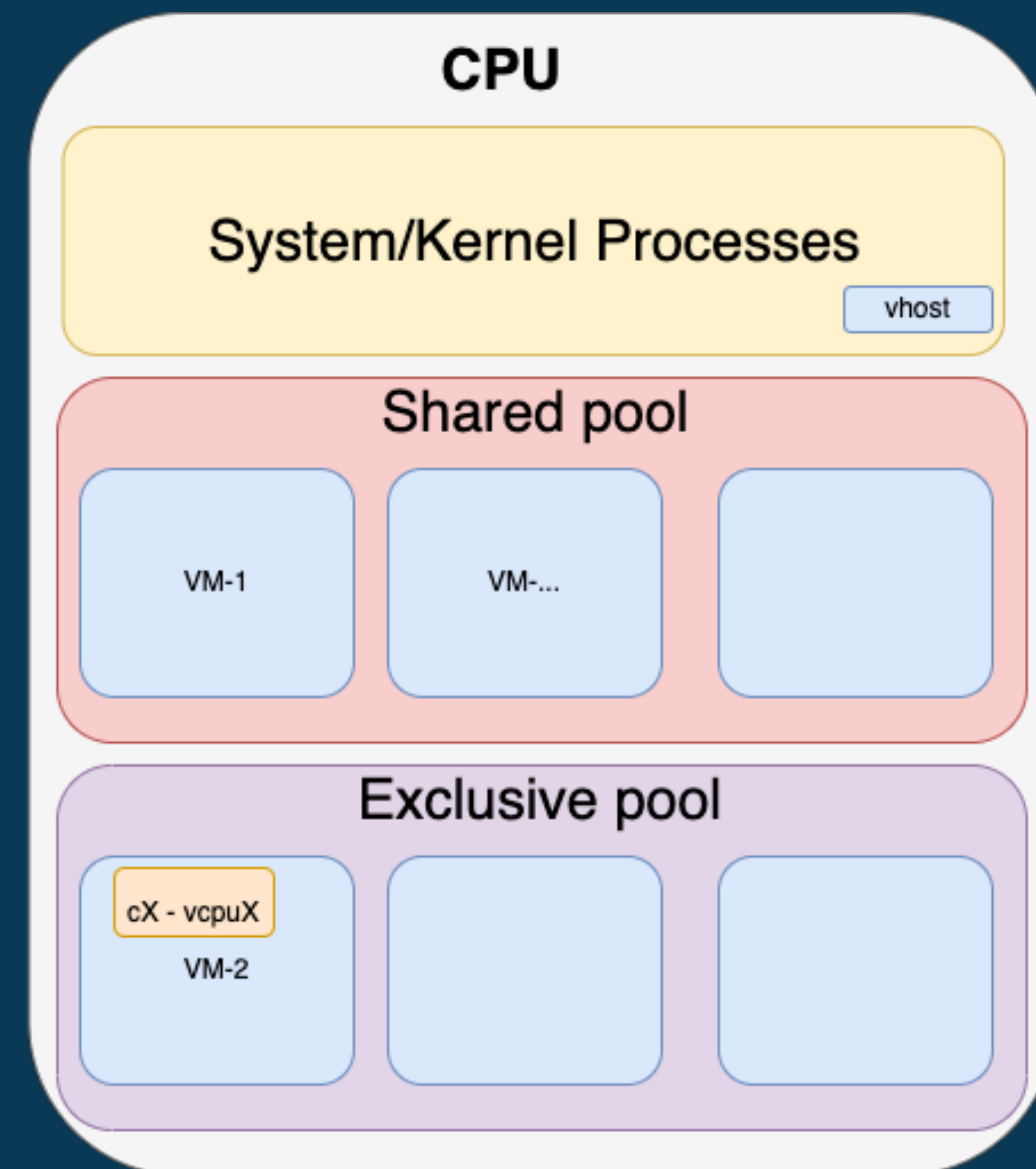
- Отдельно есть пул процессоров, которые обслуживают виртуальные машины.
- Отдельно есть пул, в котором находятся системные процессы и всякие kworker, которые обслуживают сеть/диски в виртуальных машинах.



Тюнинг CPU

Изоляция: пул для виртуальных машин поделен на два

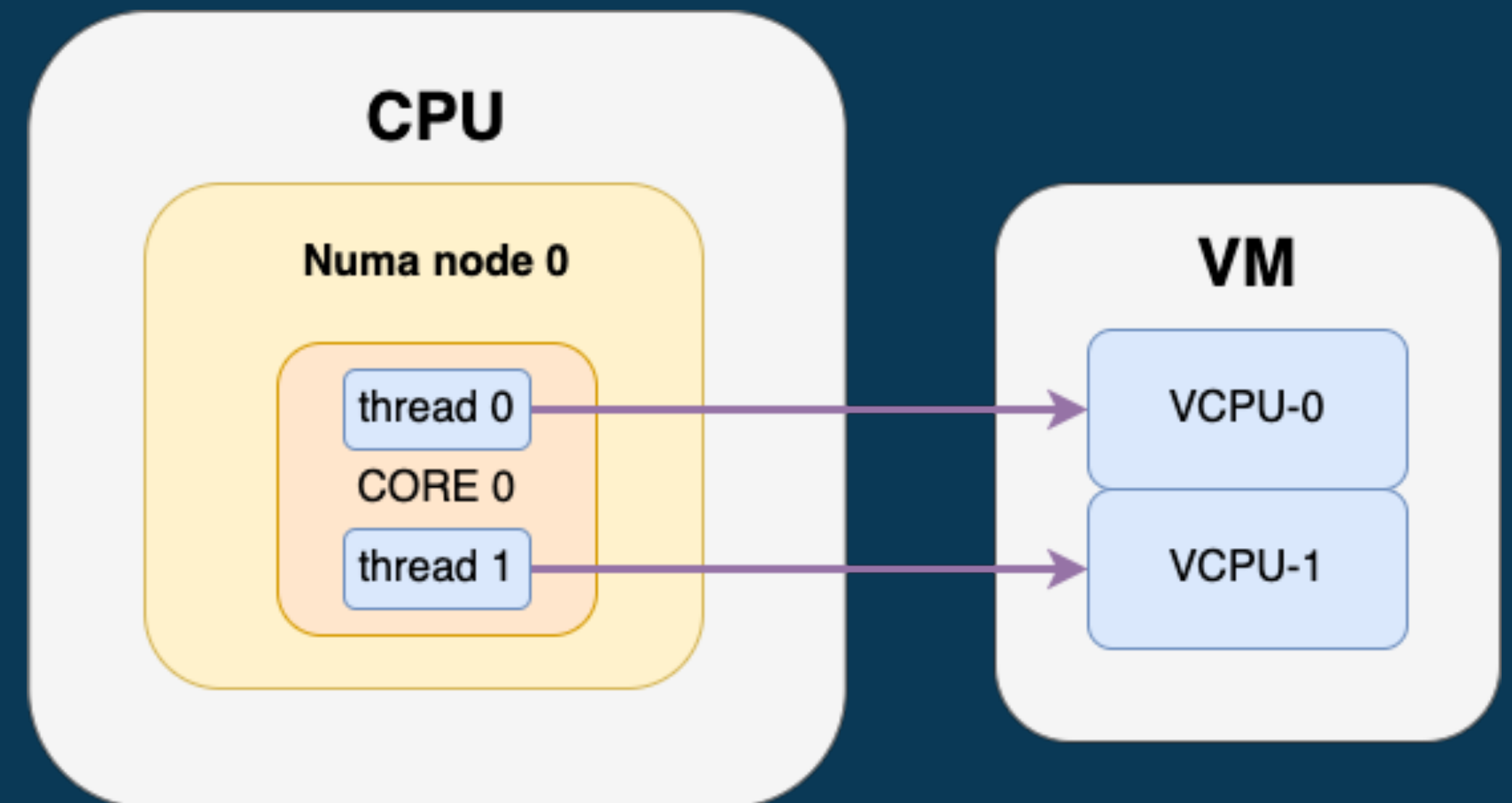
- **Эксклюзивный**
 - в нем происходит прямое, эксклюзивное выделение каждому vspu виртуальной машины к реальному CPU
 - имеет эластичный размер, если нет VM с эксклюзивными CPU, то размер пула 0
- **Шареный**
 - все остальные VM варятся в общем котле



Тюнинг CPU

Изоляция: учитываем топологию CPU

- Для системного пула
- Для эксклюзивного пула



Тюнинг CPU

Отключение SMT (Hyper-Threading): 30% по stddev pg_stat_statements

- Для shared pool — деградация производительности.
- Для exclusive pool — увеличение производительности.

```
echo off > /sys/devices/system/cpu/smt/control
```



Тюнинг CPU

Эластичный System Pool

- В системном пуле кроме vhost-net находятся emulator pin (до +20% stddev pgss для машин в шареном пуле)
- Если system + user cpu time превышает отметку 75%, то накидываем в системный пул дополнительные ядра

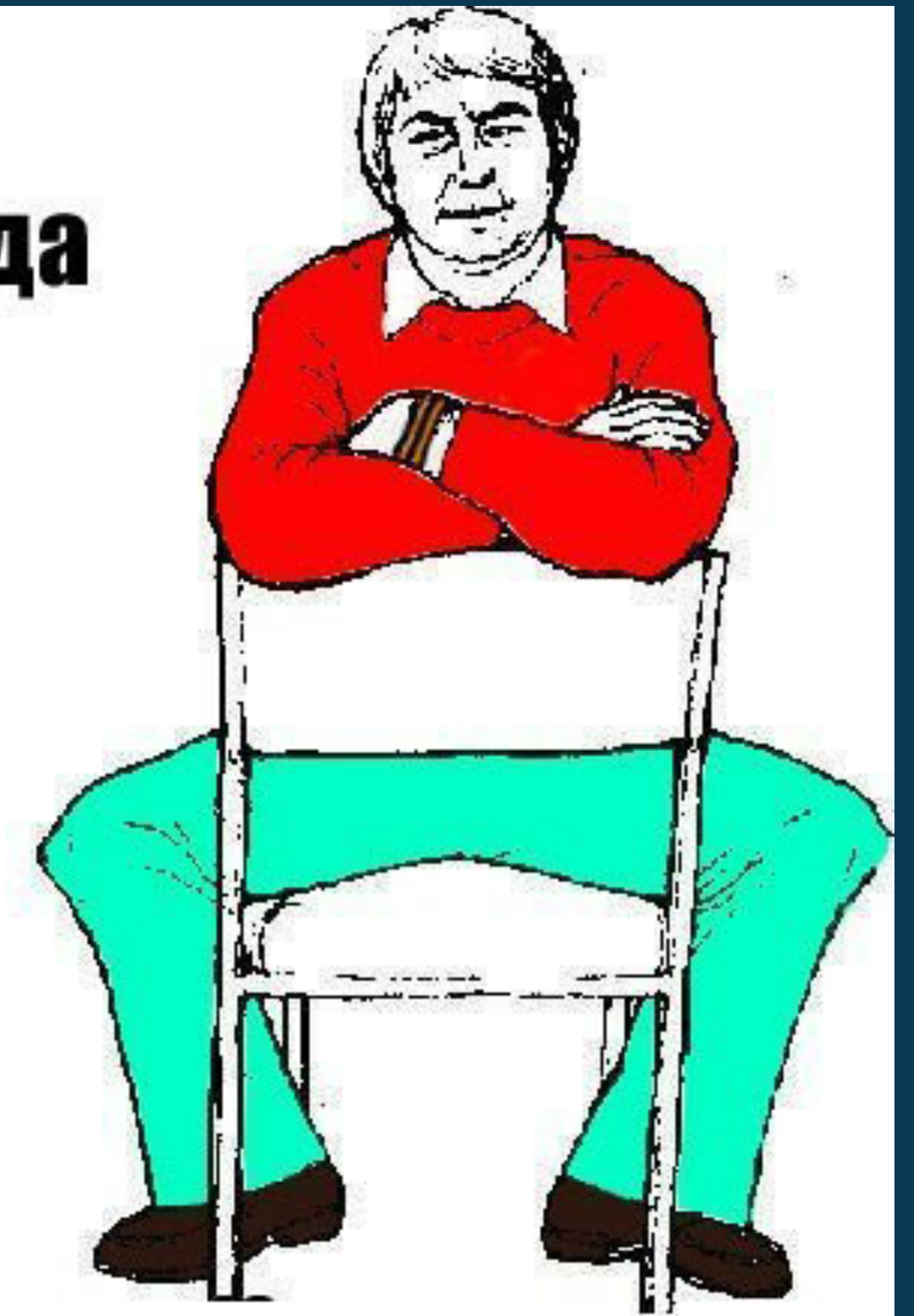
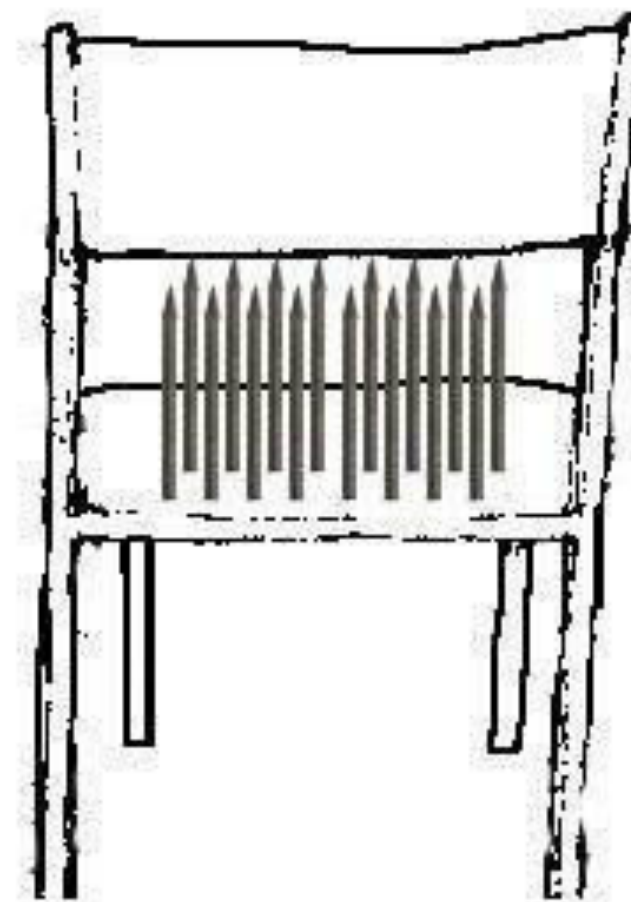


Тюнинг сети

Проблема ещё не решена

- Bridge public mode
 - случайно фиксировали ЧУЖОЙ трафик в своей VM от соседней VM, всплески systime
- Macvtap private mode
 - не видим соседей, проблема с шардированием
- SR-IOV
 - неоднородное железо
 - иногда течет libvirt
 - проблема с некоторыми Intel-карточками

Выбор есть всегда



<https://access.redhat.com/solutions/2822941>

<https://access.redhat.com/solutions/4847951>

ozon{tech

Спасибо за внимание

Васильев Дмитрий, DBA

dmitrivasilyev@ozon.ru