



# Менеджер контейнерів

## islander

Денис Герасимук, Ярослав Морозевич, Дмитро Лопушанський



# Суть проекту

Розробити аналог docker'a, який зможе запускати процеси в повністю ізольованих середовищах. **islander має такі функції:**

- Обмеження використання файлової системи, процесорної завантаженості, пам'яті, мережі
- Налаштування cgroups і створення namespace'ів
- client-server архітектура з можливістю запуску клієнта і сервера на різних хостах + шифрування каналу спілкування
- Менеджмент контейнерів, монтування директорій, volumes, підтримка cloud-стореджу та ін

# Етапи розробки



1

Етап дослідження:  
Технології, функціонал, деталі  
реалізації.

2

Розробка скриптів, які  
зможуть ізолювати процеси  
за певними параметрами

3

Написання парсера і  
сервера, які будуть  
спілкуватися через сокети

4

Поєднання всіх  
частин проекту

5

Підтримка Volumes,  
Bind Mount, TmpFS  
Менеджмент даних контейнера

6

Додавання network  
namespace, менеджмент  
контейнерів

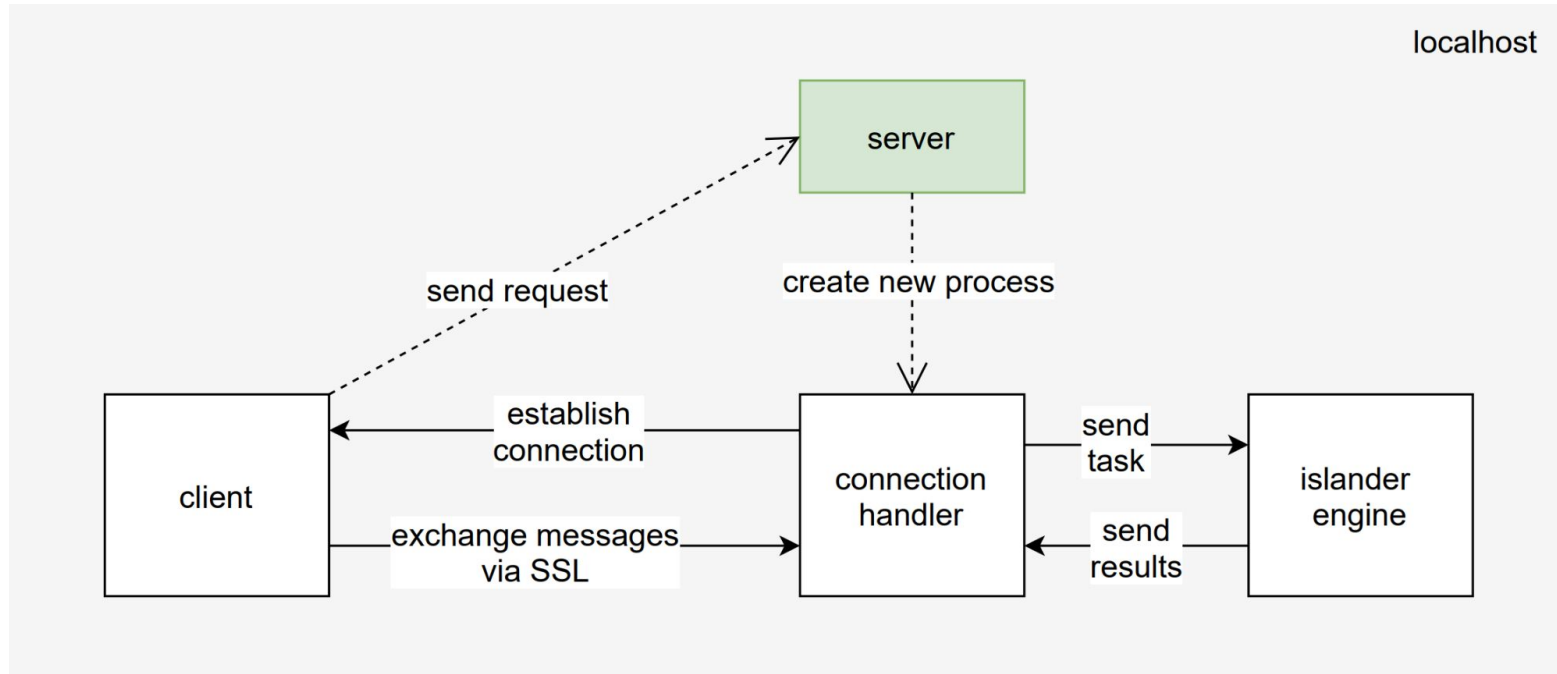
7

Покращення взаємодії  
клієнтів та сервера; робота з  
новими типами програм

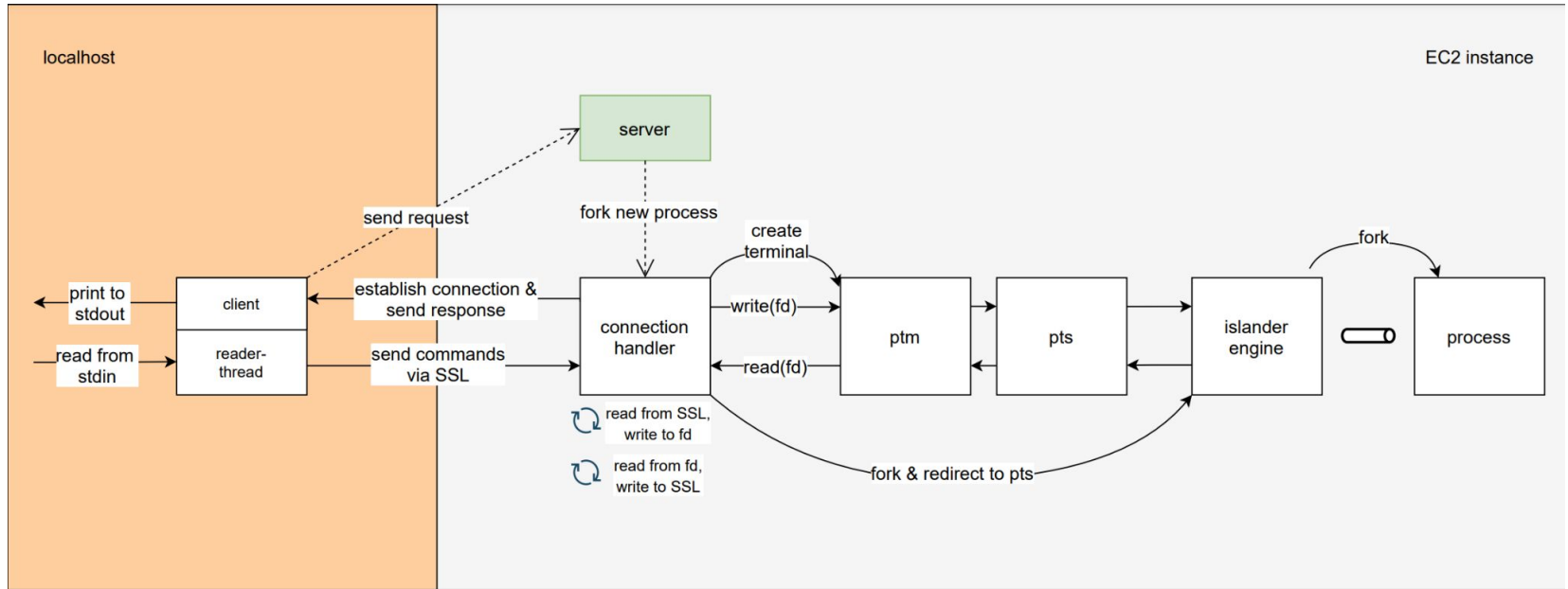
8

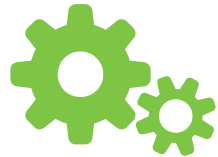
Використання хмарних  
сервісів для запуску наших  
контейнерів (EC2, S3)

# Початкова архітектура



# Архітектура проекту

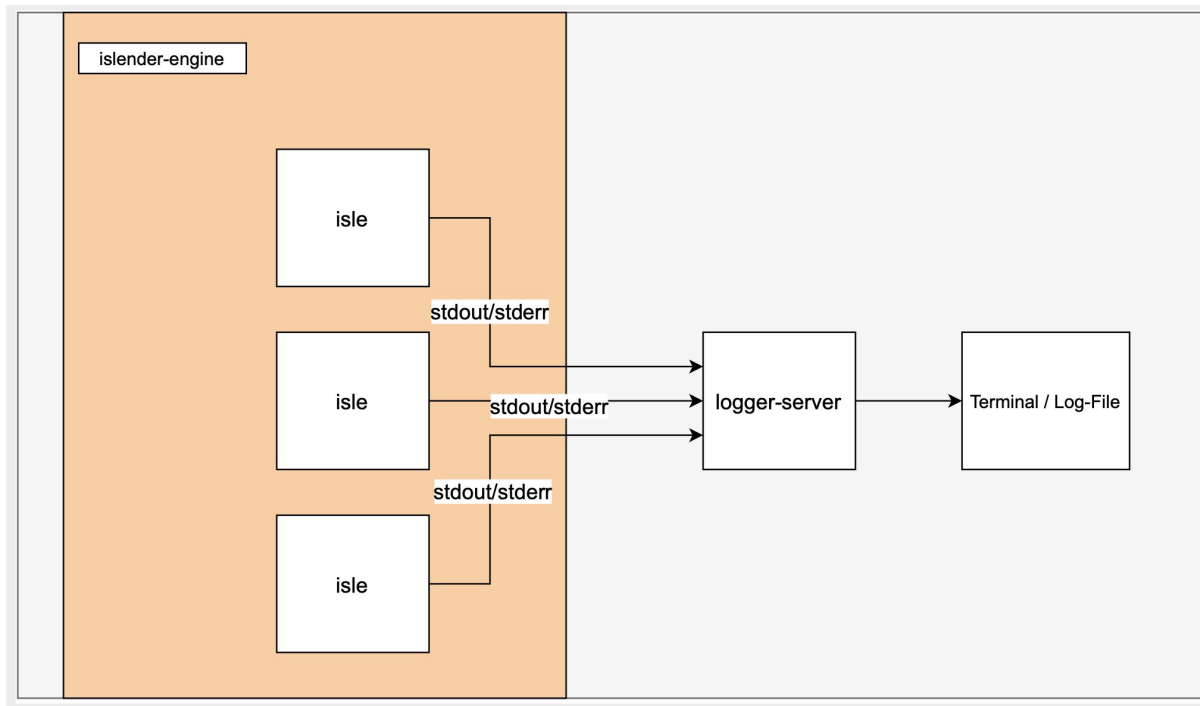
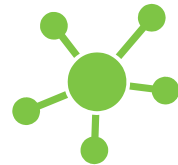


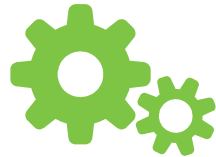


# Cgroups & Namespaces

- **Namespace** — обмежує привілеї процесу
  - Mount/Network/UTS/PID/User Namespace
- **Cgroup** — ставить ліміти та обмежує типи ресурсів
  - blkio/cpu/devices/net\_cls/memory
- **Менеджмент контейнерів**
  - ps/delete/detach
- **Менеджмент даних**
  - Volumes/BTRFS/TMPFS

# Логер-сервер





# Islander Data Management

## Проблеми:

- Дані не зберігаються, коли цього контейнера більше не існує
- Спільний доступ до даних
- high-performance I/O тощо

## Рішення:

- Islander volumes
- Bind mounts
- Tmpfs mounts



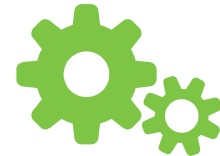
1-1 Global public cloud purchases will be \$236 billion in 2020, 23% higher than our 2014 forecast

Total public cloud revenues  
(US\$ billions)

■ Total public cloud revenues, 2016 forecast  
■ Total public cloud revenues, 2014 forecast



\*Forrester forecast

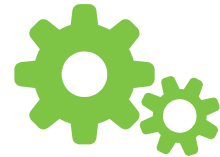


# Islander Remote Volumes

Підтримка наступних провайдерів:

- AWS S3
- Azure Storage
- GCP Cloud Storage





# Remote Volumes

**Реалізація:** Terraform + Cloud Mount Utils

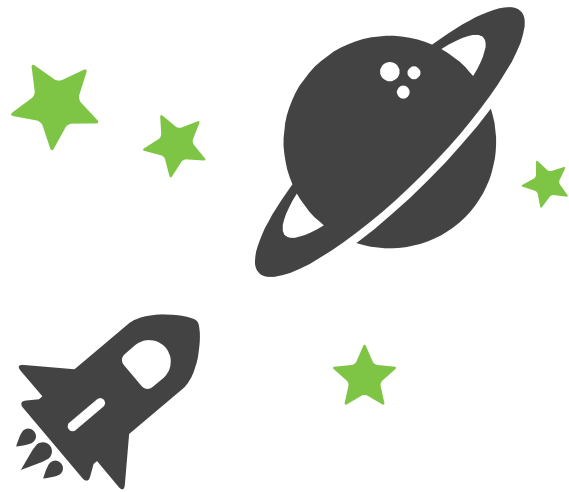
## **Terraform:**

- Infrastructure as Code
- Найбільша підтримка cloud функціоналу
- Support + documentation

Повна демонстрація remote volumes –  
<https://youtu.be/ljaCSn4vsOo>



# Демо

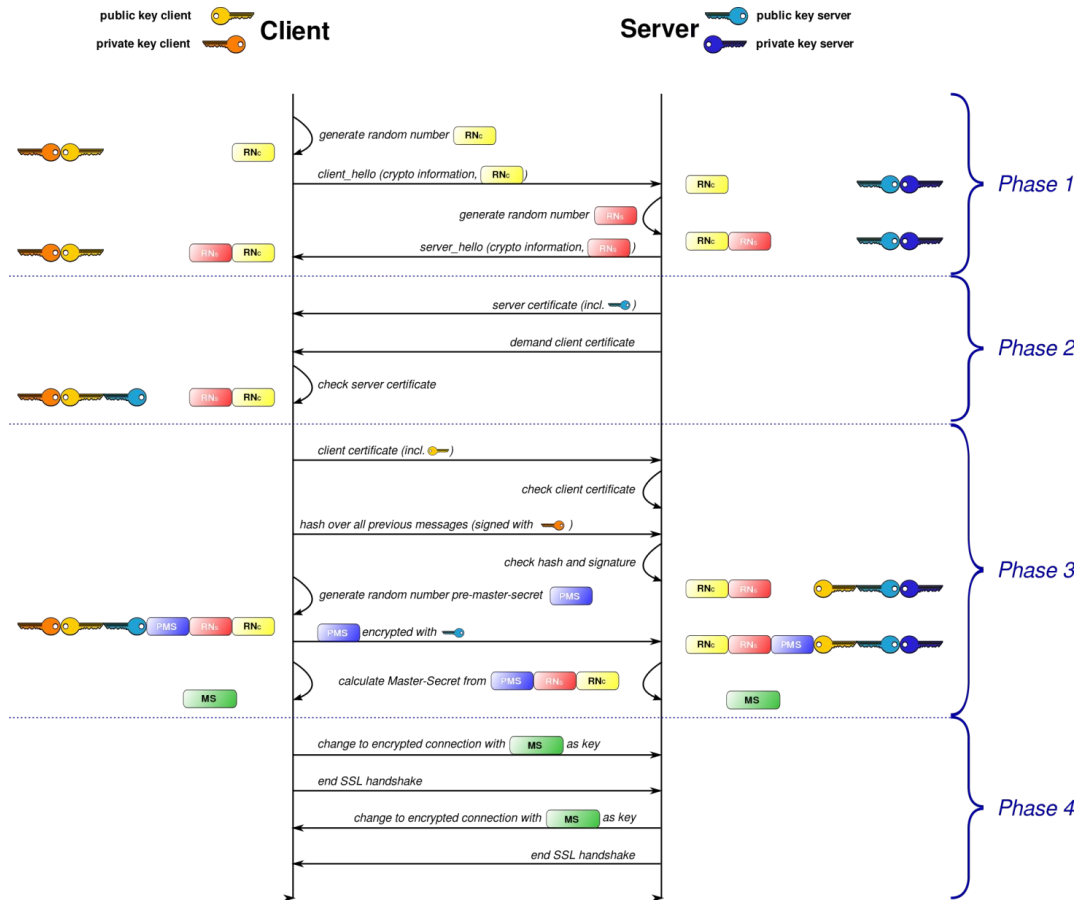


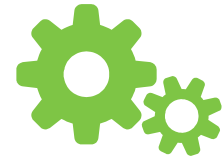
# Дякуємо!

Час для запитань

[https://github.com/denysgerasymuk799/UCU\\_OS\\_Course\\_Project](https://github.com/denysgerasymuk799/UCU_OS_Course_Project)



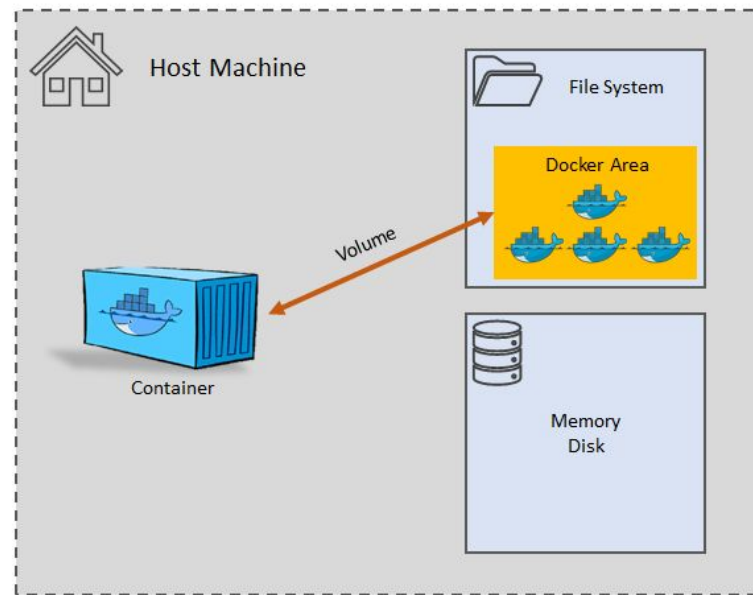


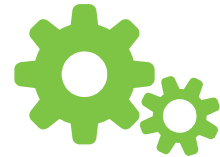


# Islander Volumes

## Особливості:

- Більша ефективність, ніж у mount
- Легше створити backup або перемістити
- Взаємодія з volumes на віддалених хостах або хмарних провайдерах





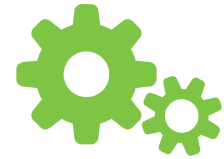
# Btrfs or B-tree

## Особливості:

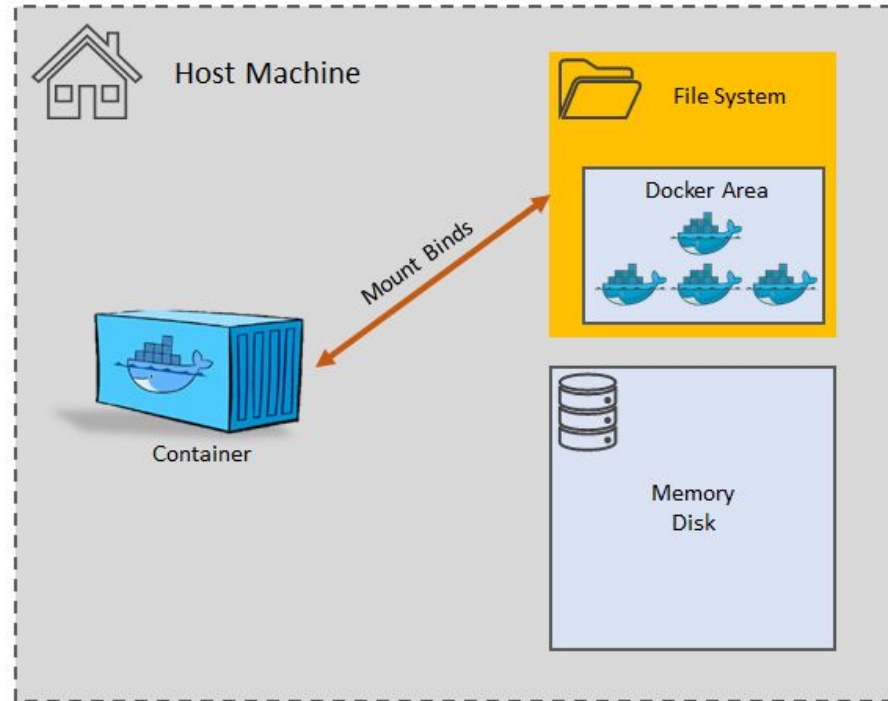
- Створення snapshots
- Використовує алгоритми компресії даних на рівні файлової системи
- Використовує контрольну суму **CRC32C** → цілісність даних і уникнути пошкодження даних
- Оптимізована підтримка SSD-дисків тощо



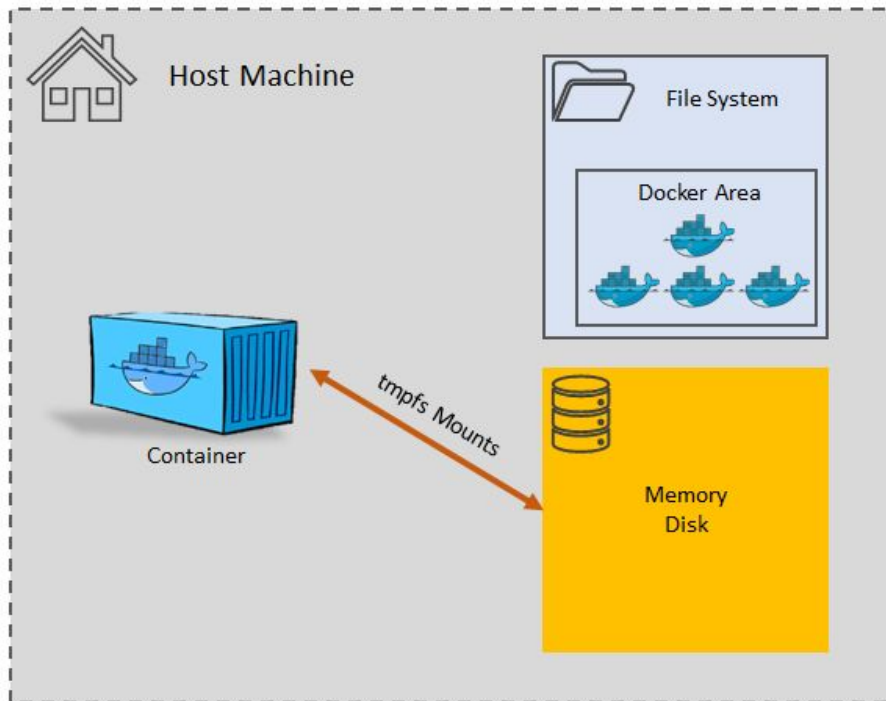
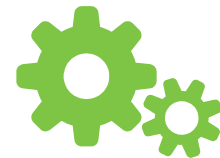


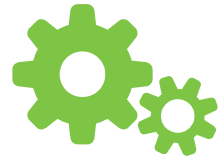


# Bind Mounts



# Tmpfs Mounts





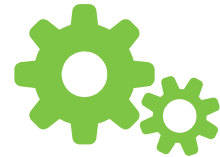
# Модель Cgroup

- Подібні до процесів
  - ієрархічні
  - дочірні cgroups успадковують певні атрибути від батьківської cgroup
- Відмінне:
  - Linux є єдиним деревом процесів
  - модель cgroup — одне або кілька окремих, не пов'язаних між собою дерев процесів



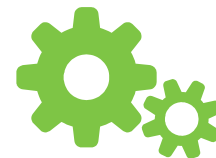
# Що таке namespace?

- Ізоляційний механізм для ресурсів
- Забезпечує відображення ресурсів зі змінами дозволів
- Зміни до процесів, які знаходяться в певному просторі імен, є невидимі поза його межами

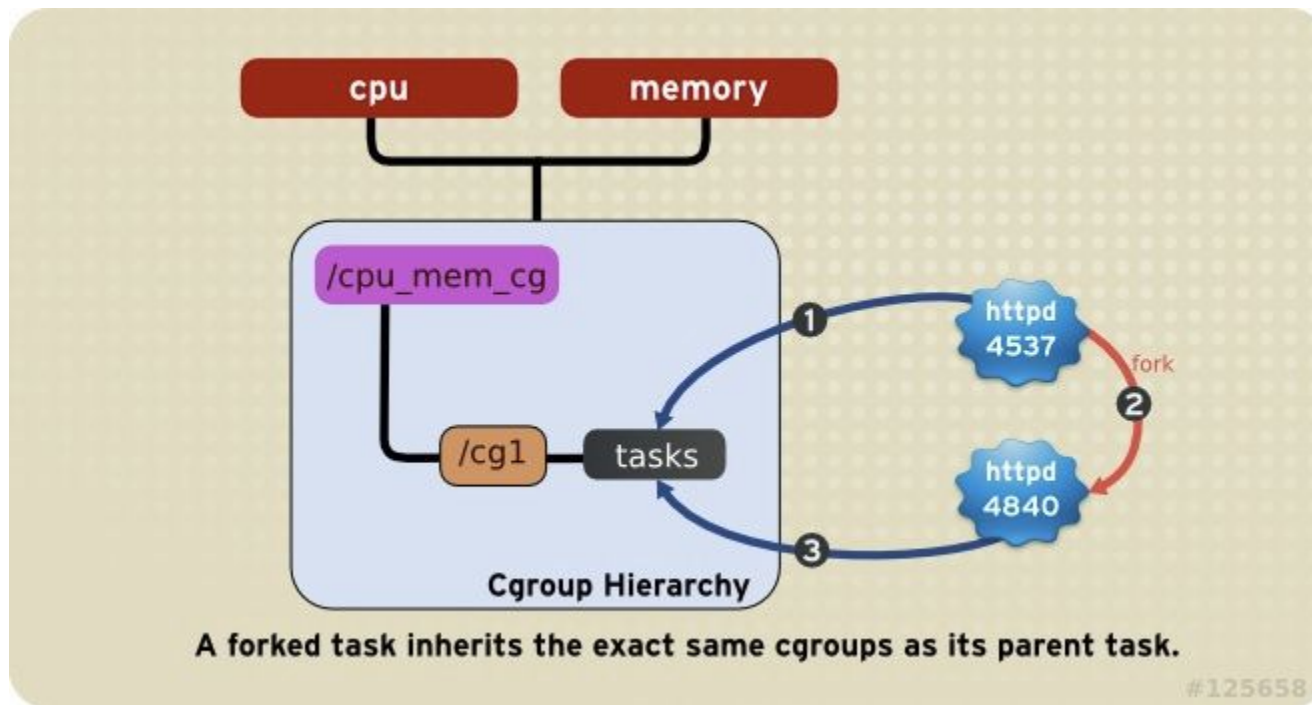


# Що таке Cgroup?

- **Namespace** — обмежує привілеї процесу  
**Cgroup** — ставить ліміти та обмежує типи ресурсів
- Дозволяють розподіляти ресурси серед визначених груп процесів

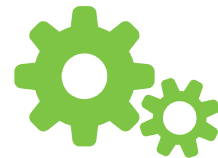


# Модель Cgroup



# Види namespace'ів

- **Mount** - керує точками монтування
- **Network** - керує мережевим стеком
- **PID** - надає процесам незалежний набір id
- **UTS** - дозволяє одній системі мати різні імена хостів/доменів
- **User Namespace** - забезпечує ізоляцію привілеїв користувача
- **IPC** - забезпечує комунікацію між процесами



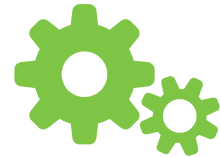
# Cgroup subsystems

- **blkio** - читання та запис блочних девайсів
- **cpu** - доступ до процесора
- **devices** - доступ до девайсів
- **net\_cls** - ліміти network io
- **memory** - RAM ліміти для cgroup

```
$ ls /sys/fs/cgroup/
```

|         |             |         |                  |            |
|---------|-------------|---------|------------------|------------|
| blkio   | cpu,cpuacct | freezer | net_cls          | perf_event |
| cpu     | cpuset      | hugetlb | net_cls,net_prio | pids       |
| cpuacct | devices     | memory  | net_prio         | systemd    |





# Приклад використання

## # Create a group

```
$ cd /sys/fs/cgroup
```

```
$ mkdir -p memory/group1
```

## # Set a memory limit of 150M

```
$ echo 150M > memory/group1/memory.limit_in_bytes
```

## # Add shell to group

```
$ echo $$ > memory/group1/tasks
```

