

**LXC**

**Linux Containers – Следующее поколение  
виртуальных машин для облака**



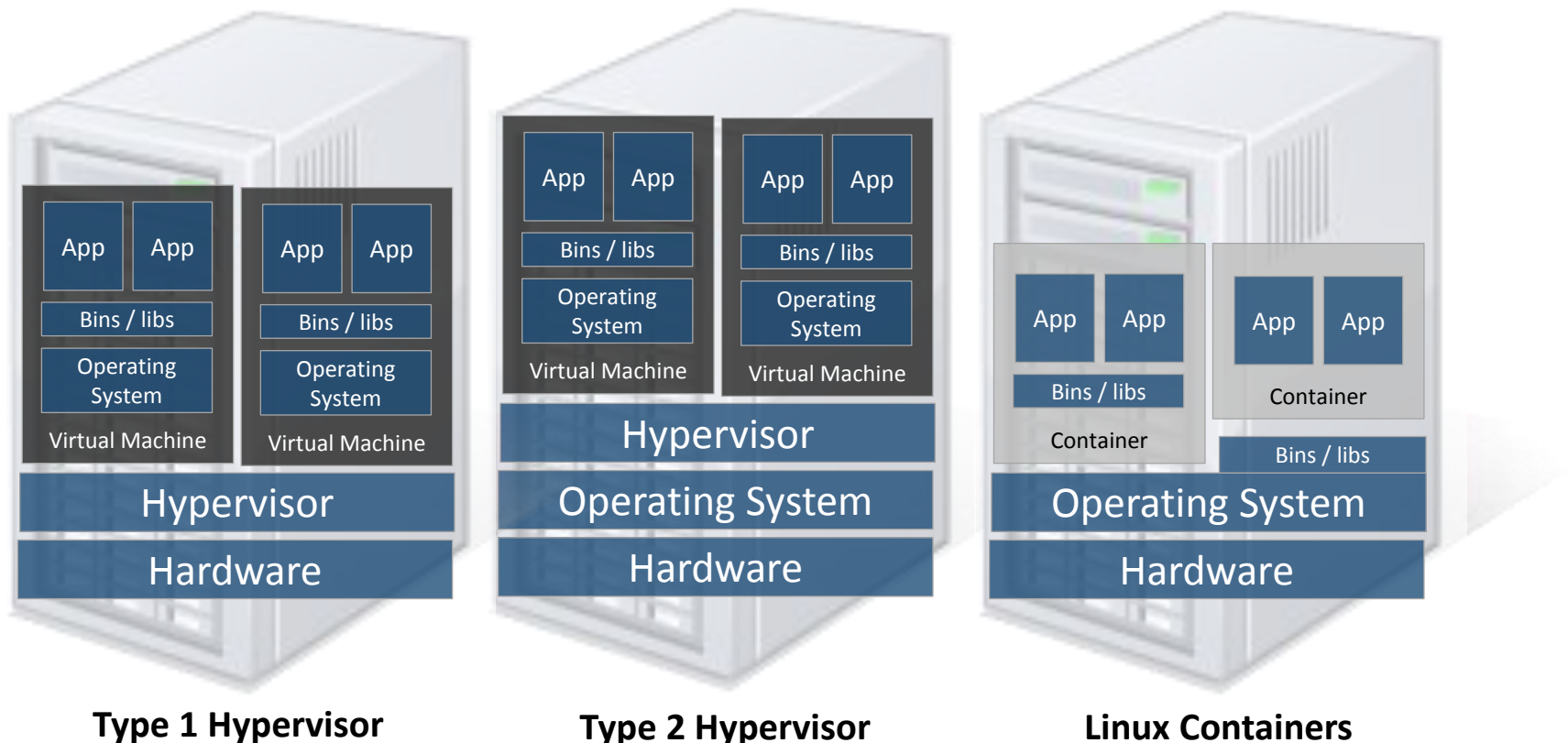
# Определения

- **Linux Containers (LXC → LinuX Containers)**
  - Легковесные виртуальные машины
  - Реализовано средствами современного ядра Linux
  - VMs без гипервизора
  
- ***Контейнер для***
  - (Linux) Операционной системы
  - Одного или нескольких приложений
  
- **LXC как технология ≠ LXC “tools”**

# Гипервизоры против Linux контейнеров

Контейнеры разделяют ресурсы хостовой операционной системы, следовательно Они легковесные, но имеют одинаковую версию ядра.

Контейнеры изолированные, но разделяют ядро и требуемые модули хостовой ОС

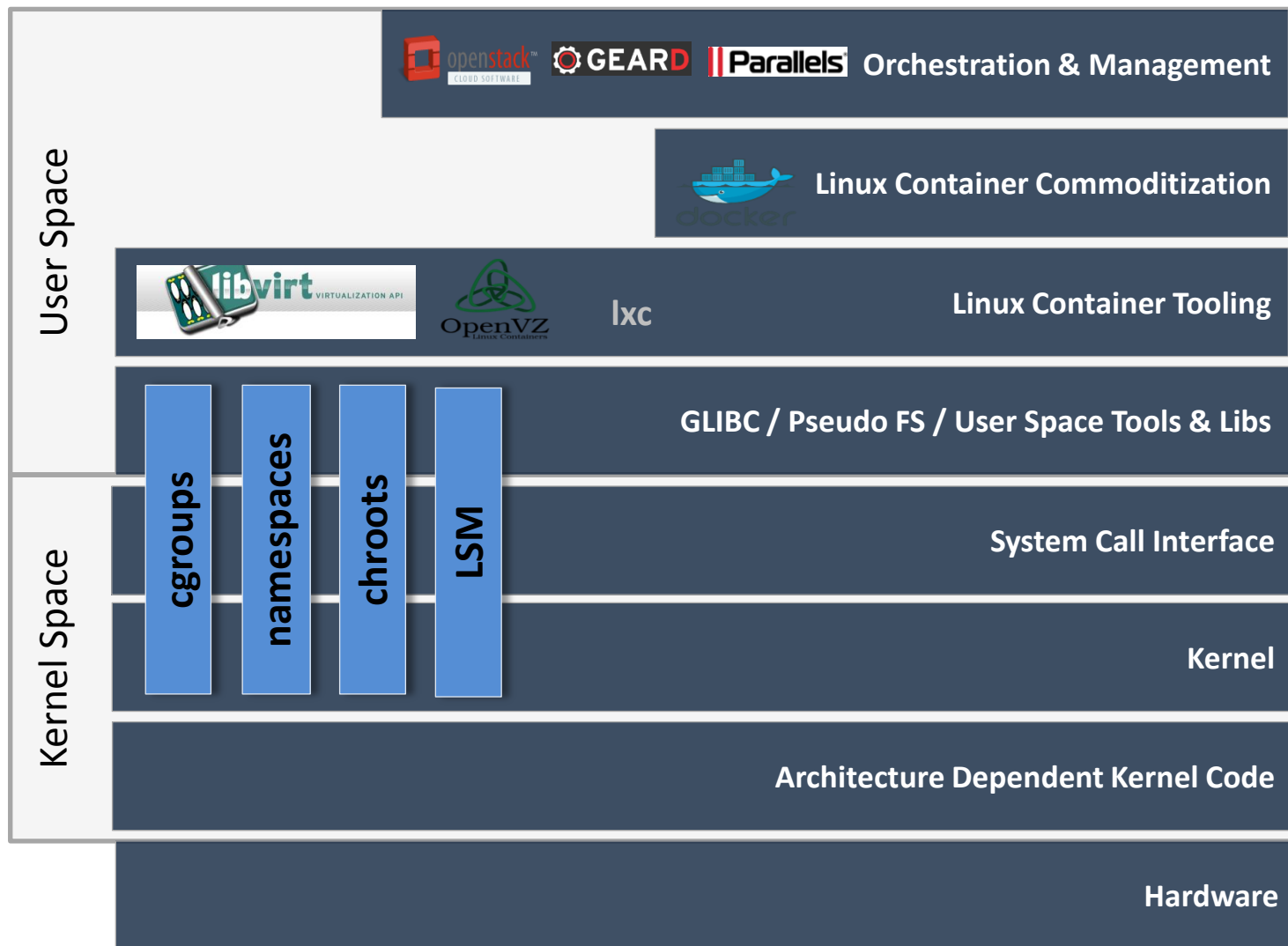


Type 1 Hypervisor

Type 2 Hypervisor

Linux Containers

# Стек технологий LXC



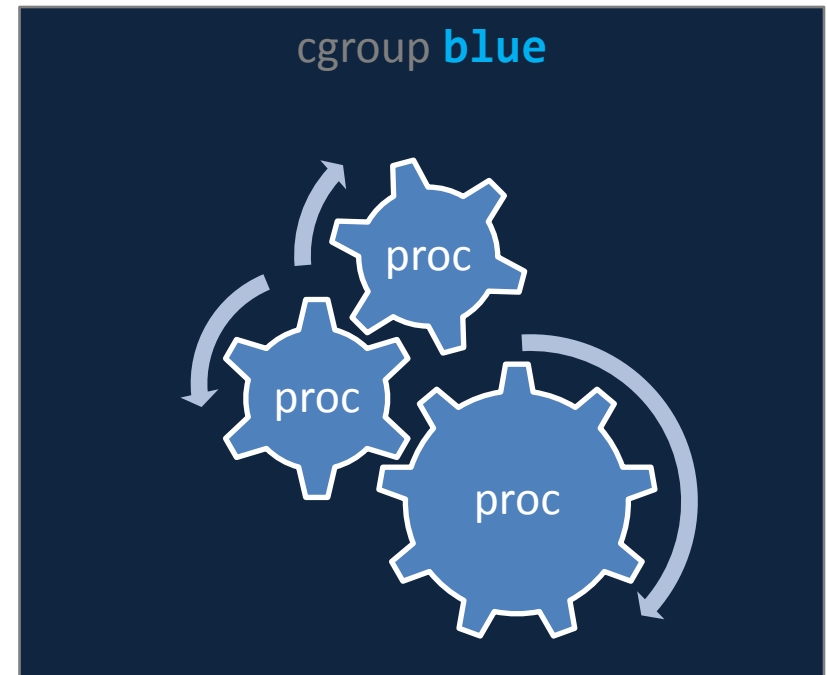
# Группы управления Linux(cgroups)

## ■ Проблема

- Как можно контролировать и анализировать выполнение задач(процессов)?

## ■ Решение → группы управления (cgroups)

- Доступ к устройствам
- Ограничения ресурсов
- Приоритеты
- Привязка
- Управление
- Вложенности



# Подсистема Linux cgroup

Подсистема	Контролируемые параметры
<b>blkio</b>	<ul style="list-style-type: none"><li>- Ограничение доступа к устройствам ввода/вывода. Группа на устройство.</li><li>- Ограничения для блочных устройств как на количество IO операций, так и на количество IO операций за секунду.</li></ul>
<b>cpu</b>	<ul style="list-style-type: none"><li>- Лимитирует процессорное время(микросекунды или секунды) для группы процессов.</li><li>- Лимит на время на одном CPU.</li><li>- Получение процессорного времени пропорционально приоритету.</li></ul>
<b>cpuset</b>	<ul style="list-style-type: none"><li>- Доступ к группе CPUs.</li><li>- Возможность мигрировать между областями ОЗУ.</li><li>- Блокировки доступа к памяти, приоритеты и тд.</li></ul>
<b>devices</b>	<ul style="list-style-type: none"><li>- Определяет к каким устройствам есть доступ.</li></ul>
<b>freezer</b>	<ul style="list-style-type: none"><li>- Приостанавливает/возобновляет работы группы процессов.</li></ul>
<b>memory</b>	<ul style="list-style-type: none"><li>- Лимит памяти для группы(в байтах).</li><li>- Подкачка, контроль памяти и тд.</li></ul>
<b>hugetlb</b>	<ul style="list-style-type: none"><li>- Лимит использования hugeTLB страниц.</li><li>- Контроль использования hugeTLB группой.</li></ul>
<b>net_cls</b>	<ul style="list-style-type: none"><li>- Разделение пакетов по id сети.</li><li>- Приоритет для пакетов группы.</li></ul>
<b>net_prio</b>	<ul style="list-style-type: none"><li>- Приоритет на сетевую активность на интерфейс.</li></ul>

# Интерфейс псевдофайловой системы Linux cgroups

- Псевдо фс Linux - интерфейс для cgroups
  - Директория на каждую подсистему cgroup
  - Чтение / запись в файл в директории своей группы

/sys/fs/cgroup/my-lxc

```
-- blkio
|  -- blkio.io_merged
|  -- blkio.io_queued
|  -- blkio.io_service_bytes
|  -- blkio.io_serviced
|  -- blkio.io_service_time
|  -- blkio.io_wait_time
|  -- blkio.reset_stats
|  -- blkio.sectors
|  -- blkio.throttle.io_service_bytes
|  -- blkio.throttle.io_serviced
|  -- blkio.throttle.read_bps_device
|  -- blkio.throttle.read_iops_device
|  -- blkio.throttle.write_bps_device
|  -- blkio.throttle.write_iops_device
|  -- blkio.time
|  -- blkio.weight
|  -- blkio.weight_device
|  -- cgroup.clone_children
|  -- cgroup.event_control
|  -- cgroup.procs
|  -- notify_on_release
|  -- release_agent
|  -- tasks
-- cpu
|  -- ...
-- ...
-- perf_event
```

echo "8:16 1048576" >  
blkio.throttle.read\_bps\_device

cat blkio.weight\_device  
dev weight  
8:1 200  
8:16 500



# Структура ФС Linux cgroups

/sys/fs/cgroup/

## -- blkio

```
-- blkio.io_merged
-- blkio.io_queued
-- blkio.io_service_bytes
-- blkio.io_serviced
-- blkio.io_service_time
-- blkio.io_wait_time
-- blkio.reset_stats
-- blkio.sectors
-- blkio.throttle.io_service_bytes
-- blkio.throttle.io_serviced
-- blkio.throttle.read_bps_device
-- blkio.throttle.read_iops_device
-- blkio.throttle.write_bps_device
-- blkio.throttle.write_iops_device
-- blkio.time
-- blkio.weight
-- blkio.weight_device
-- cgroup.clone_children
-- cgroup.event_control
-- cgroup.procs
-- lxc
-- notify_on_release
-- release_agent
-- tasks
```

## -- cpu

```
-- cgroup.clone_children
-- cgroup.event_control
-- cgroup.procs
-- cpu.cfs_period_us
-- cpu.cfs_quota_us
-- cpu.rt_period_us
-- cpu.rt_runtime_us
-- cpu.shares
-- cpu.stat
-- lxc
-- notify_on_release
-- release_agent
-- tasks
```

## -- cpuacct

```
-- cgroup.clone_children
-- cgroup.event_control
-- cgroup.procs
-- cpuacct.stat
-- cpuacct.usage
-- cpuacct.usage_percpu
-- lxc
-- notify_on_release
-- release_agent
-- tasks
```

## -- cpuset

```
-- cgroup.clone_children
-- cgroup.event_control
-- cgroup.procs
-- cpuset.cpu_exclusive
-- cpuset.cpus
-- cpuset.mem_exclusive
-- cpuset.mem_hardwall
-- cpuset.memory_migrate
-- cpuset.memory_pressure
-- cpuset.memory_pressure_enabled
-- cpuset.memory_spread_page
-- cpuset.memory_spread_slab
-- cpuset.mems
-- cpuset.sched_load_balance
-- cpuset.sched_relax_domain_level
-- lxc
-- notify_on_release
-- release_agent
-- tasks
```

## -- devices

```
-- cgroup.clone_children
-- cgroup.event_control
-- cgroup.procs
-- devices.allow
-- devices.deny
-- devices.list
-- lxc
-- notify_on_release
-- release_agent
-- tasks
```

## -- freezer

```
-- cgroup.clone_children
-- cgroup.event_control
-- cgroup.procs
-- lxc
-- notify_on_release
-- release_agent
-- tasks
```

## -- hugetlb

```
-- cgroup.clone_children
-- cgroup.event_control
-- cgroup.procs
-- hugetlb.2MB.failcnt
-- hugetlb.2MB.limit_in_bytes
-- hugetlb.2MB.max_usage_in_bytes
-- hugetlb.2MB.usage_in_bytes
-- lxc
-- notify_on_release
-- release_agent
-- tasks
```

## -- memory

```
-- cgroup.clone_children
-- cgroup.event_control
-- cgroup.procs
-- lxc
-- memory.failcnt
-- memory.force_empty
-- memory.limit_in_bytes
-- memory.max_usage_in_bytes
-- memory.memsw.failcnt
-- memory.memsw.limit_in_bytes
-- memory.memsw.max_usage_in_bytes
-- memory.memsw.usage_in_bytes
-- memory.move_charge_at_immigrate
-- memory.numa_stat
-- memory.oom_control
-- memory.soft_limit_in_bytes
-- memory.stat
-- memory.swappiness
-- memory.usage_in_bytes
-- memory.use_hierarchy
-- notify_on_release
-- release_agent
-- tasks
```

## ■ cgroup subsystem mounts

- Root of cgroup subsystem for a given cgroup controller
- All cgroup instances are children of root
- All processes initially in tasks

## ■ cgroup instance

- Individual cgroup instance
- Contains a full set of pseudo files for the subsystem
- Hierarchical controls influenced by parent
- Individual processes added to tasks pseudo file

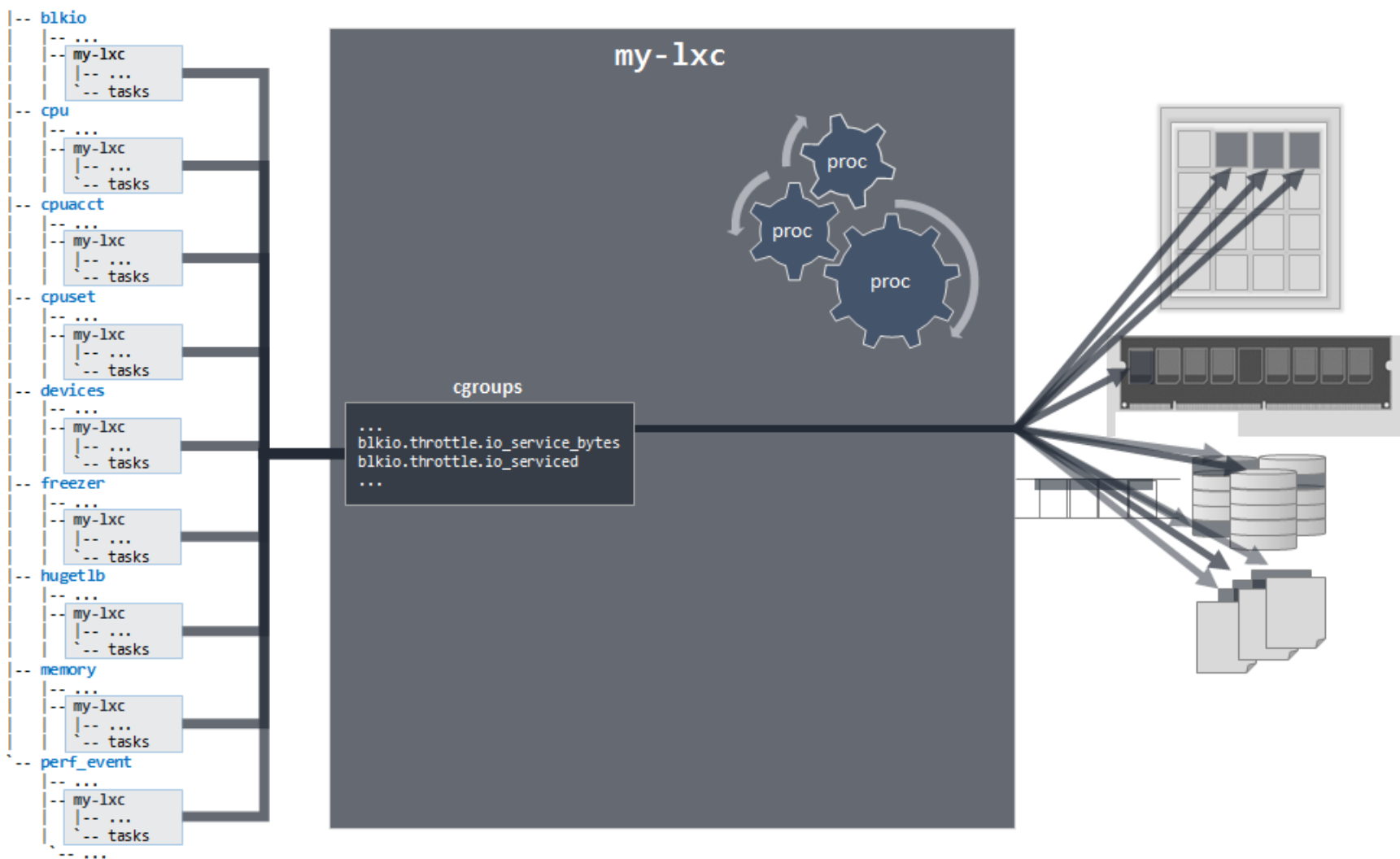
## User space management

- Cgroup interface is the pseudo FS itself
- Read counter files to obtain metrics (cat, head, etc.)
- Write values into control files

## -- perf\_event

```
-- cgroup.clone_children
-- cgroup.event_control
-- cgroup.procs
-- lxc
-- notify_on_release
-- release_agent
-- tasks
```

# Как строится контейнер?



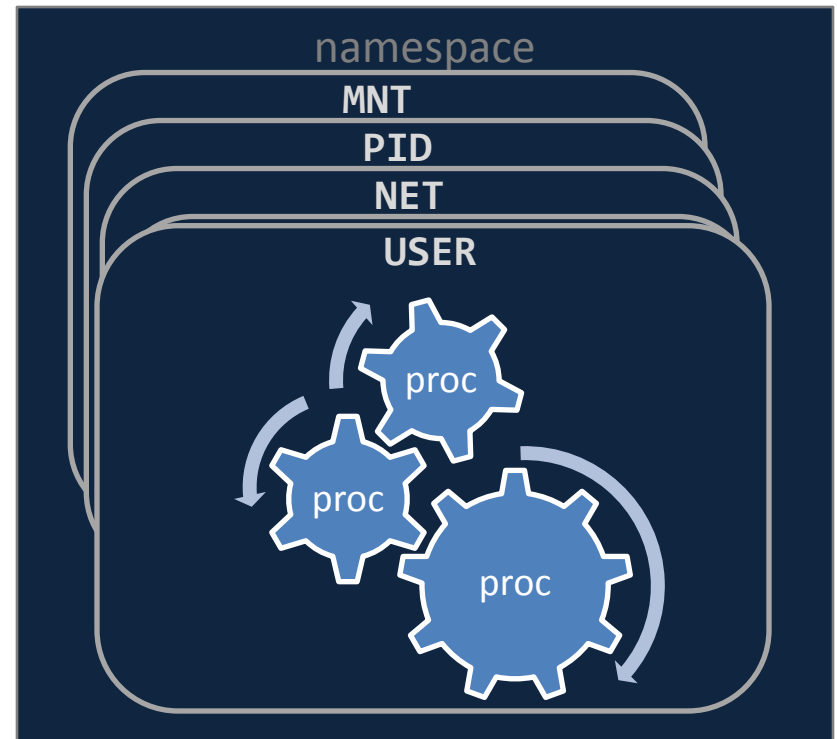
# Пространства имен

## ■ Проблема

- Как предоставить изолированное представление глобальных ресурсов задачам(процессам)

## ■ Решение → пространства имен

- **MNT**; точки монтирования, файловые системы и тд.
- **PID**; процессы
- **NET**; маршрутизация и тд..
- **IPC**; System V IPC
- **UTS**; хостовое и доменное имя
- **USER**; UID и GID



# Пространства имен Linux - концепция

## global (i.e. root) namespace

### MNT NS

```
/
/proc
/mnt/fsrd
/mnt/fsrw
/mnt/cdrom
/run2
```

### UTS NS

```
globalhost
rootns.com
```

### PID NS

PID	COMMAND
1	/sbin/init
2	[kthreadd]
3	[ksoftirqd]
4	[cpuset]
5	/sbin/udev
6	/bin/sh
7	/bin/bash

### IPC NS

SHMID	OWNER
32452	root
43321	boden

SEMID	OWNER
0	root
1	Boden

MSQID	OWNER
-------	-------

### NET NS

```
lo: UNKNOWN...
eth0: UP...
eth1: UP...
br0: UP...

app1 IP:5000
app2 IP:6000
app3 IP:7000
```

### USER NS

```
root 0:0
ntp 104:109
mysql 105:110
boden 106:111
```

## purple namespace

### MNT NS

```
/
/proc
/mnt/purplenfs
/mnt/fsrw
/mnt/cdrom
```

### UTS NS

```
purplehost
purplens.com
```

### PID NS

PID	COMMAND
1	/bin/bash
2	/bin/vim

### IPC NS

SHMID	OWNER

SEMID	OWNER
0	root

MSQID	OWNER
-------	-------

### NET NS

```
lo: UNKNOWN...
eth0: UP...

app1 IP:1000
app2 IP:7000
```

### USER NS

```
root 0:0
app 106:111
```

## blue namespace

### MNT NS

```
/
/proc
/mnt/cdrom
/bluens
```

### UTS NS

```
bluehost
bluens.com
```

### PID NS

PID	COMMAND
1	/bin/bash
2	python
3	node

### IPC NS

SHMID	OWNER

SEMID	OWNER

MSQID	OWNER
-------	-------

### NET NS

```
lo: UNKNOWN...
eth0: DOWN...
eth1: UP

app1 IP:7000
app2 IP:9000
```

### USER NS

```
root 0:0
app 104:109
```

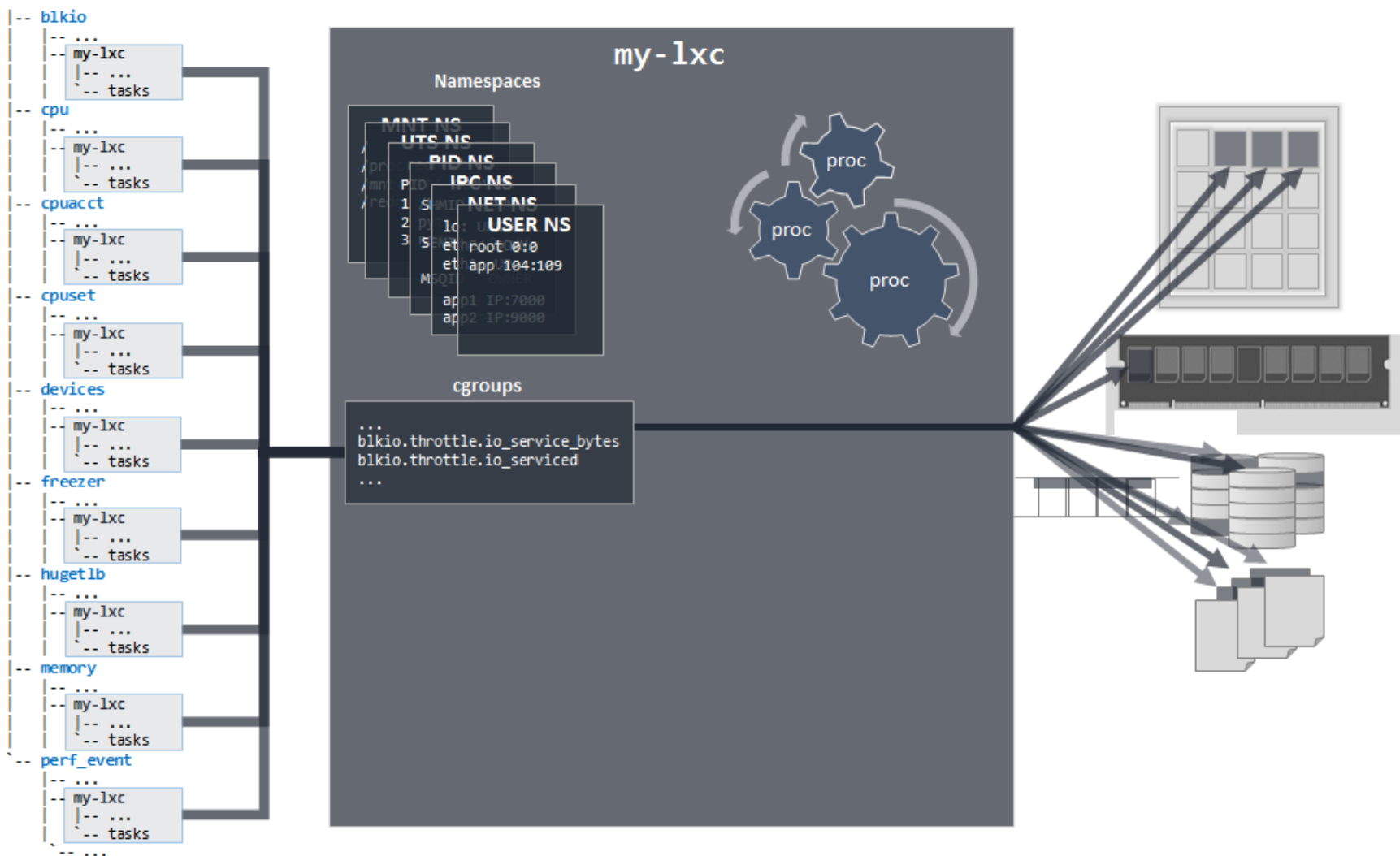
# Linux пространства имен & cgroups: Доступность

	Kernel 2.4.19	Kernel 2.6.19	Kernel 2.6.24	Kernel 2.6.29	Kernel 3.3	Kernel 3.8	Kernel 3.10+
namespaces							
Mount	✓	✓	✓	✓	✓	✓	✓
UTS		✓	✓	✓	✓	✓	✓
IPC		✓	✓	✓	✓	✓	✓
PID			✓	✓	✓	✓	✓
Network				✓	✓	✓	✓
User						✓	✓
cgroups							
Initial support			✓	✓	✓	✓	✓
net_prio					✓	✓	✓
net_cls					✓	✓	✓
Blkio async I/O							✓

Distribution	Kernel	Notes
Ubuntu 12.04.02 (LTS)	3.5.0-23.35	
Ubuntu 13.10	3.11	User namespace disabled
RHEL 6.4	2.6.32	
RHEL 7 (2014)	3.7	
Debian 7	3.2	User namespace enabled but not fully supported
Gentoo 20122112	3.6.8	
Fedora 18	3.6.10	
Fedora 19	3.10.x	User namespace enabled but not fully supported
Linux Mint 13 (LTS)	3.2.14	
Mageia 3 (STS)	3.8.13	

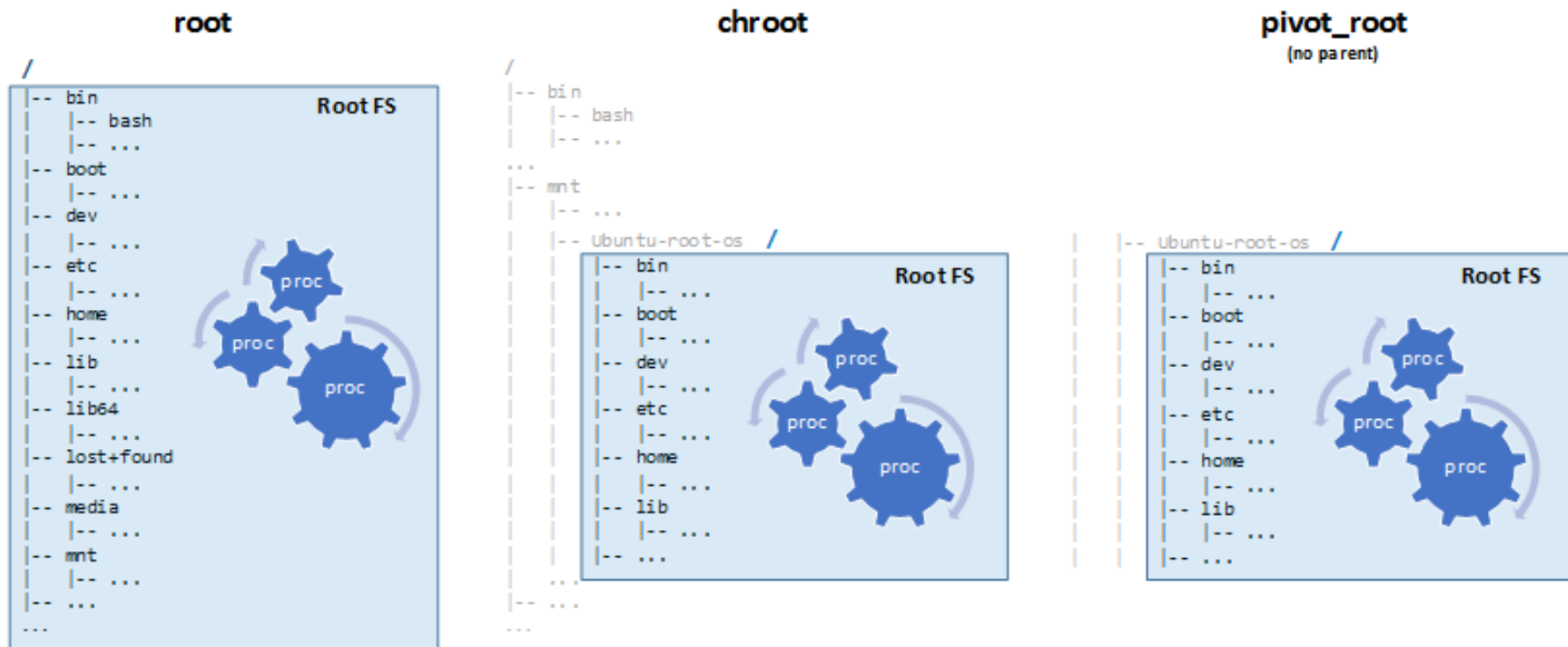
- Заметка: Пользовательские пространства имен поддерживаются в ядре 3.8+, но дистрибутивы не поддерживают:
- Отражение LXC UID/GID между контейнером и хостом
  - Non-root LXC создание

# Как строится контейнер?

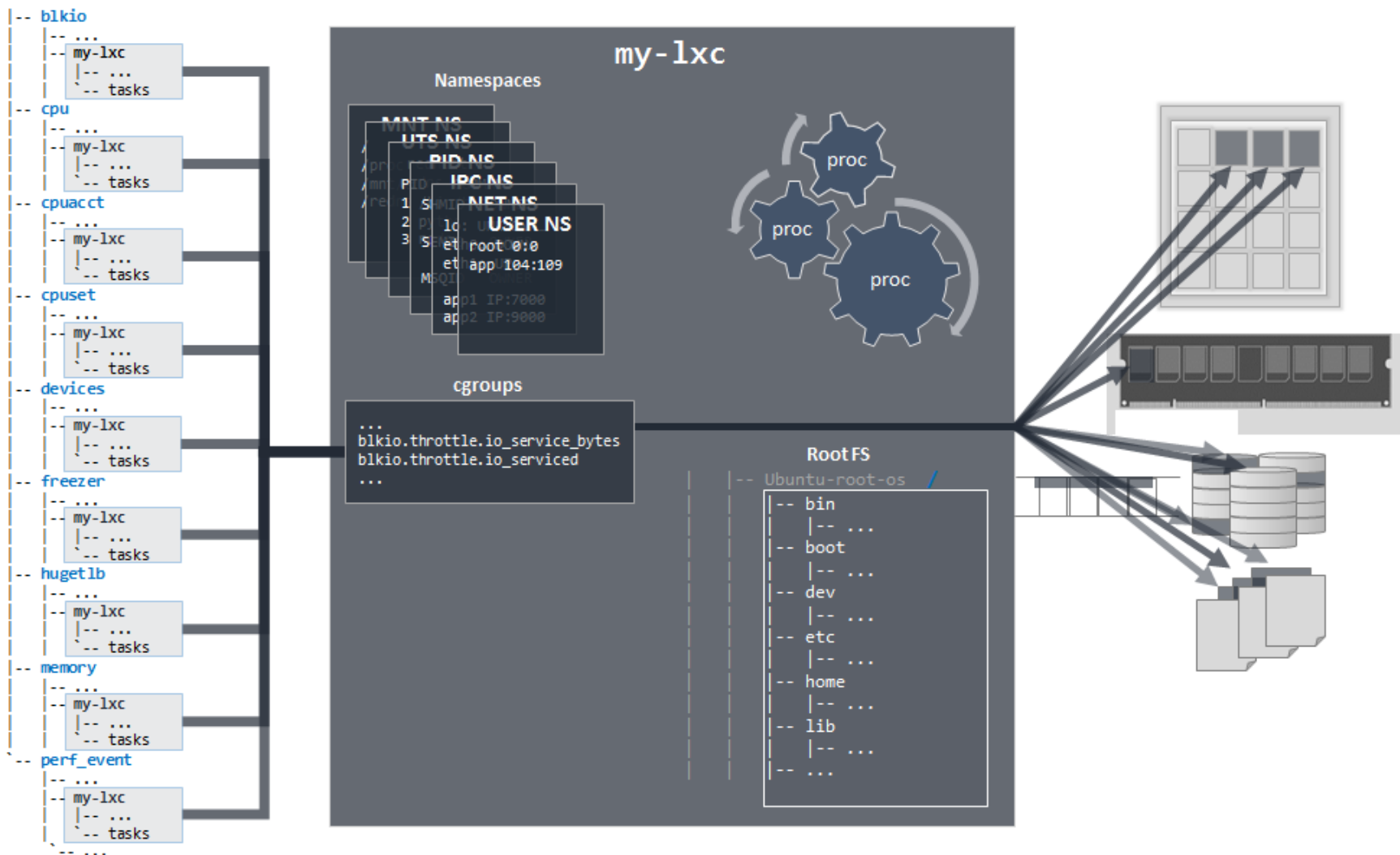


# Linux chroot & pivot\_root

- Использование `pivot_root` перекрывает использование `chroot`
- Точка привязки `pivot_root` становится новой ФС

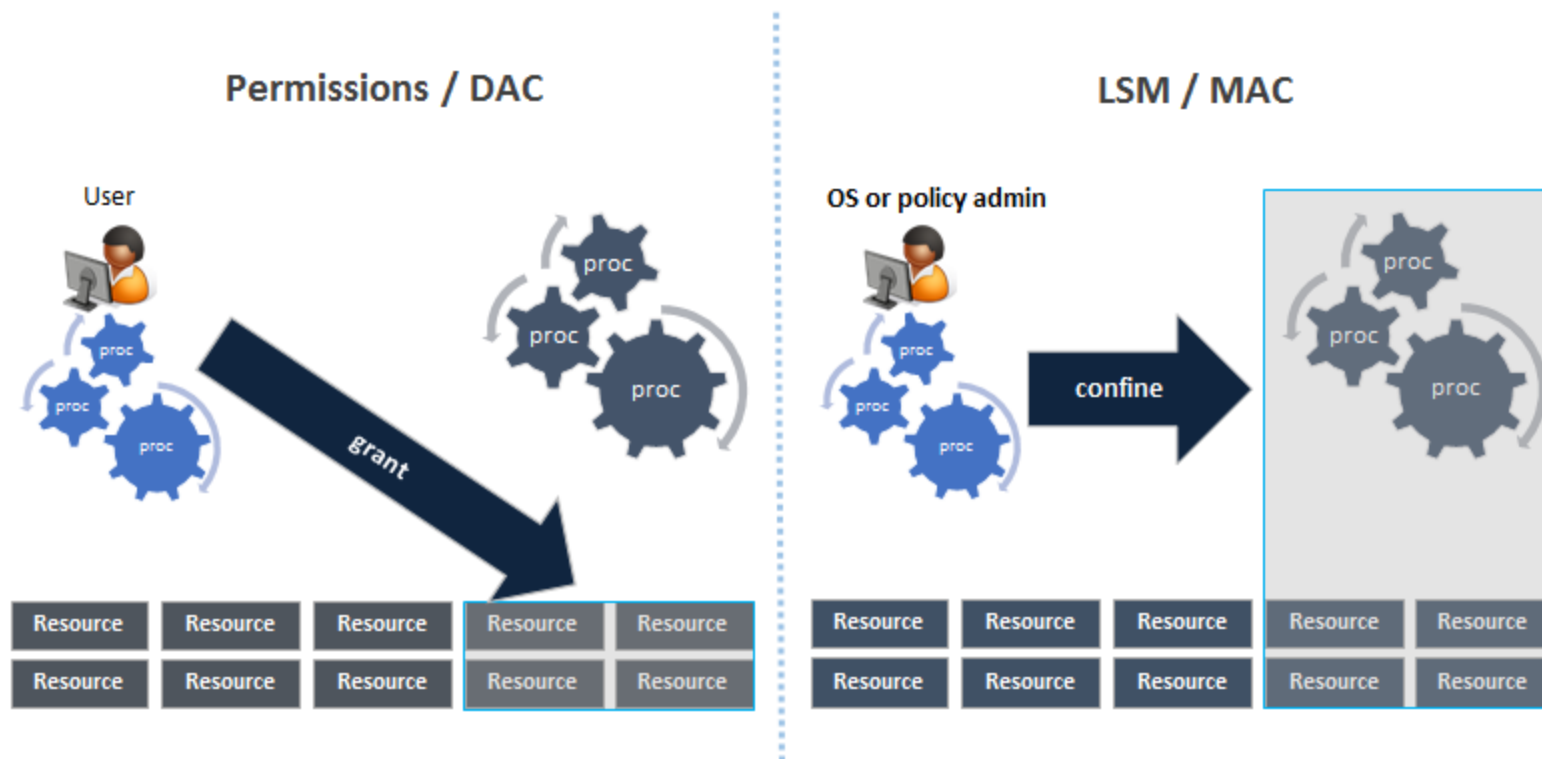


# Как строится контейнер?



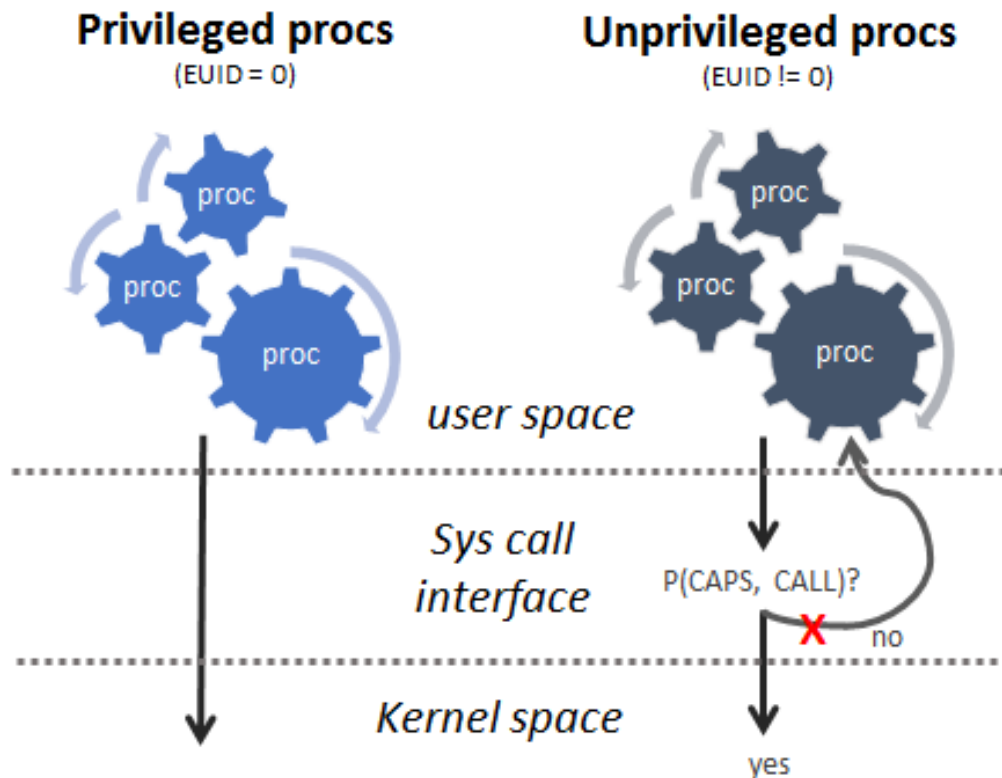
# Безопасность системы Linux & MAC

- Linux Security Modules (LSM) – модули ядра, предоставляющие MAC
- MAC vs DAC
  - В MAC, Администратор(пользователь или процесс) предоставляет доступ объекту/инициатору
  - В DAC, Владелец ресурса предоставляет полномочия процессу
- Существует поддержка LSM, включая: AppArmor, SELinux, GRSEC, etc.



# Возможности Linux

- Привилегии процесса определяются системными вызовами
- Может быть привязан к LXC процессу(ам)



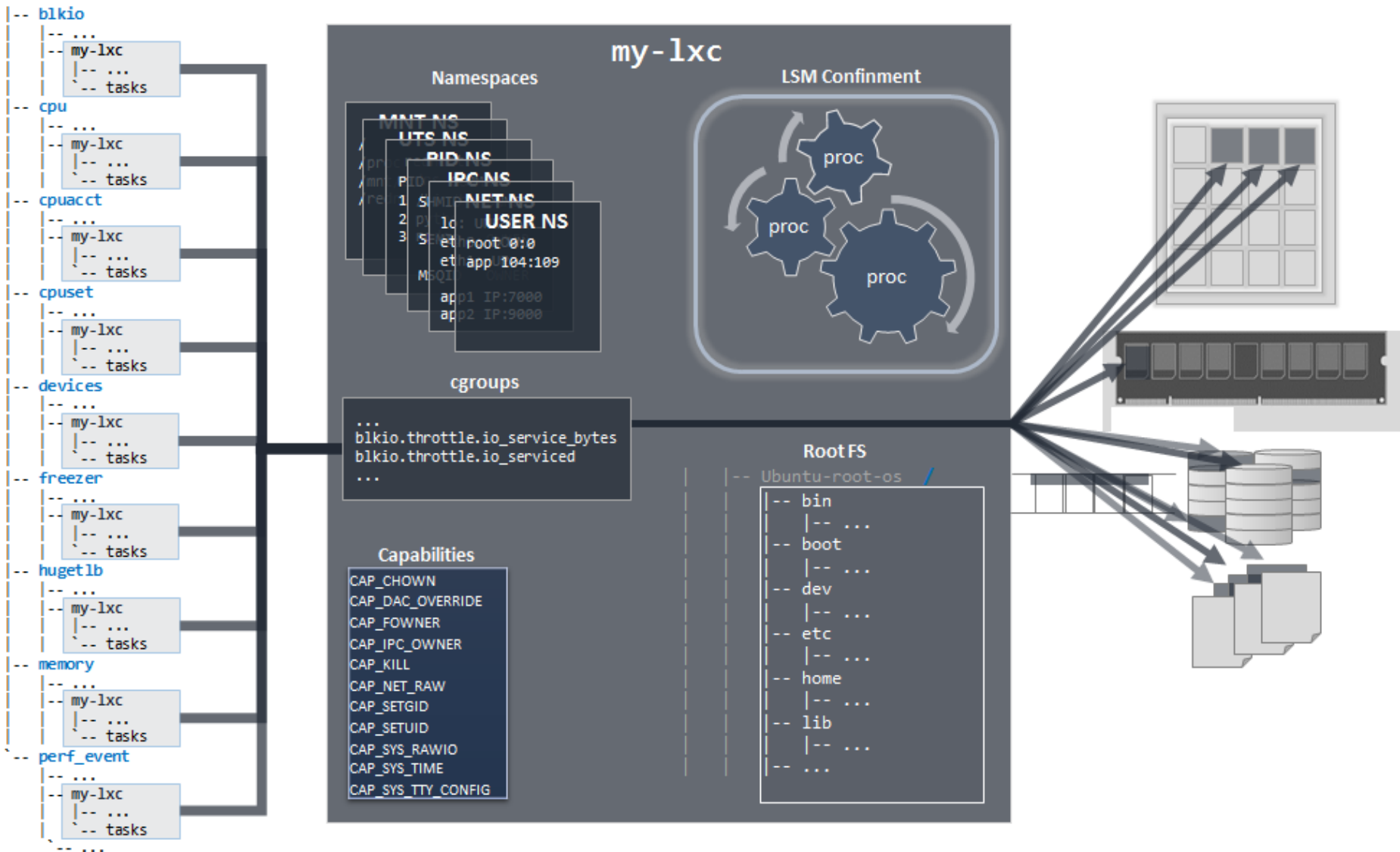
## Capabilities

```
CAP_AUDIT_CONTROL
CAP_AUDIT_WRITE
CAP_CHOWN
CAP_DAC_OVERRIDE
CAP_DAC_READ_SEARCH
CAP_FOWNER
CAP_FSETID
CAP_IPC_LOCK
CAP_IPC_OWNER
CAP_KILL
CAP_LEASE
CAP_LINUX_IMMUTABLE
CAP_MKNOD
CAP_NET_ADMIN
CAP_NET_BIND_SERVICE
CAP_NET_BROADCAST
CAP_NET_RAW
CAP_SETGID
CAP_SETPCAP
CAP_SETUID
CAP_SYS_ADMIN
CAP_SYS_BOOT
CAP_SYS_CHROOT
CAP_SYS_MODULE
CAP_SYS_NICE
CAP_SYS_PACCT
CAP_SYS_PTRACE
CAP_SYS_RAWIO
CAP_SYS_RESOURCE
CAP_SYS_TIME
CAP_SYS_TTY_CONFIG
```

# Другие меры безопасности

- Сокращение общего доступа к ФС с использованием привязок монтирования
- Linux seccomp
- Ограничить системные вызовы
- Обновлять ядро Linux
- Пользовательские пространства имен в 3.8+ ядрах
- Запуск контейнеров в качестве пользователя без полномочий root
- Отображение UID / GID в контейнере

# Как строится контейнер?



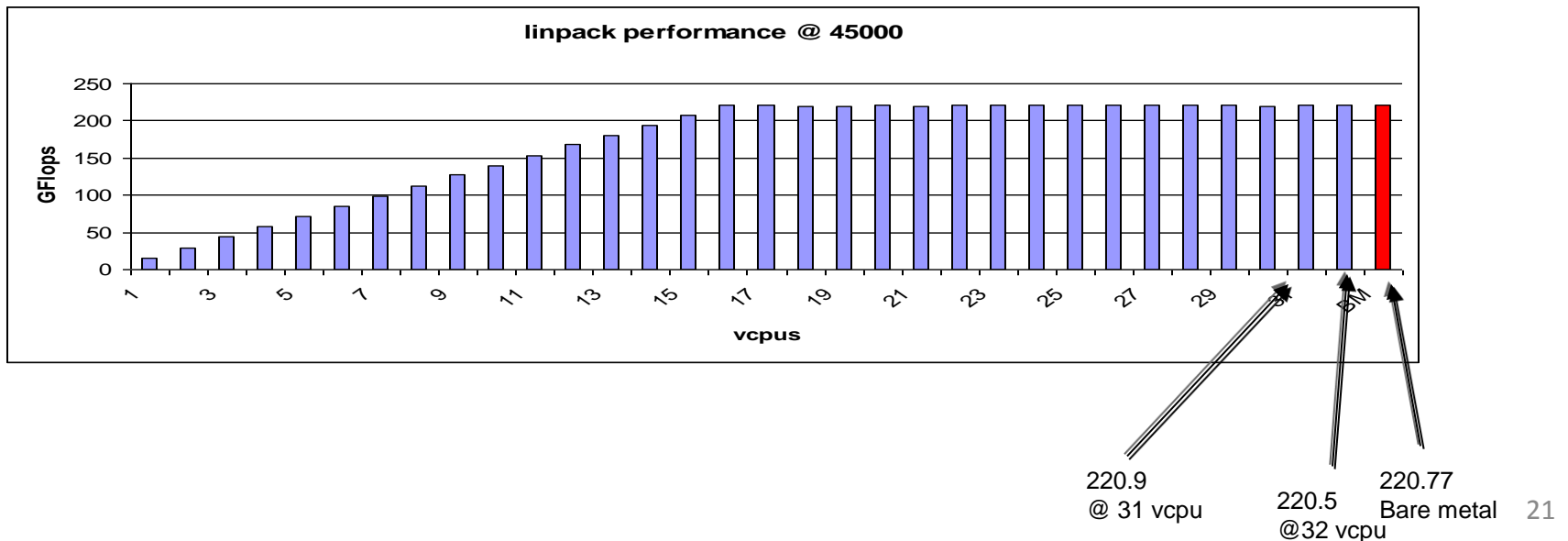
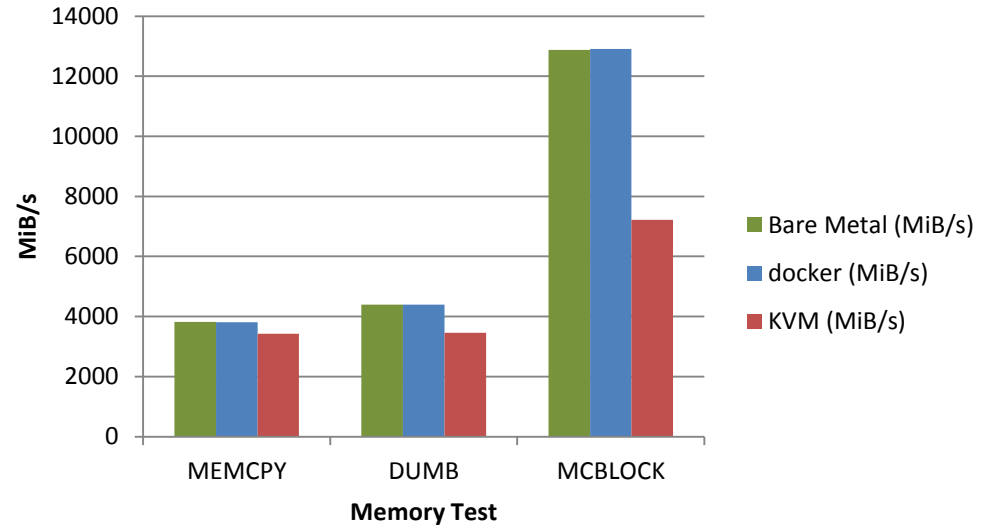
# LXC Проприетарные инструменты

	Virtuozzo	OpenVZ	Linux VServer	Libvirt-lxc	Lxc (tools)	Warden	Imctfy	Docker
Summary	Commercial product using OpenVZ under the hood	Custom Kernel providing well seasoned LXC support	A set of kernel patches providing LXC. Not based on cgroups or namespaces.	Libvirt support for LXC via cgroups and namespaces.	Lib + set of user spaces tools /bindings for LXC.	LXC management tooling used by CF.	Similar to LXC, but provides more intent based focus.	Commoditization of LXC adding support for images, build files, etc.
Part of upstream Kernel?	No	No	Partial	Yes	Yes	Yes	Yes, but additional patches needed for specific features.	Yes
License	Commercial	GNU GPL v2	GNU GPL v2	GNU LGPL	GNU LGPL	Apache v2	Apache v2	Apache v2
APIs / Bindings		<ul style="list-style-type: none"> <li>- CLI</li> <li>- API</li> </ul>	<ul style="list-style-type: none"> <li>- CLI</li> <li>- C</li> </ul>	<ul style="list-style-type: none"> <li>- CLI</li> <li>- C</li> <li>- Python</li> <li>- Java</li> <li>- C#</li> <li>- PHP</li> </ul>	<ul style="list-style-type: none"> <li>- Python</li> <li>- Lua</li> <li>- GO</li> <li>- CLI</li> </ul>			<ul style="list-style-type: none"> <li>- GO</li> <li>- REST</li> <li>- CLI</li> <li>- Python</li> <li>- Other 3<sup>rd</sup> party libs</li> </ul>
Management plane/ Dashboard	Virtuozzo Parrallels	Virtuozzo Parrallels + others		<ul style="list-style-type: none"> <li>- OpenStack</li> <li>- Archipel</li> <li>- Virt-Manager</li> </ul>	<ul style="list-style-type: none"> <li>- LXC web panel</li> <li>- Lexy</li> </ul>			<ul style="list-style-type: none"> <li>- OpenStack</li> <li>- Shipyard</li> <li>- Docker UI</li> </ul>

# Гостевая ОС: как на железном сервере

- Обычный LXC контейнер дает производительность как и реальный сервер

### Тестирование использования памяти



# LXC Итоги

- Производительность как и на голом железе
- Быстрые вычисления в облаке
- Снижение потребления ресурсов (CPU, MEM) на вычислительном узле

# LXC недостатки

*Есть несколько недочетов...*

- Отсутствие инструментов и поддержки для корпораций
- Живая миграция по-прежнему остается незавершенной
- Полная оркестровка между ресурсами (вычисление / хранение / сеть)
- Опасения по поводу безопасности
- Не известная технология ... пока
- Интеграция с существующей виртуализацией и облачным оснащением
- Не так много стандартов
- Отсутствует набор навыков
- И т.п.

# Использованная литература

- <http://www.slideshare.net/BodenRussell/realizing-linux-containerslxc>
- <http://bodenr.blogspot.com/2014/05/kvm-and-docker-lxc-benchmarking-with.html>
- <https://www.docker.io/>
- <http://sysbench.sourceforge.net/>
- <http://dag.wiee.rs/home-made/dstat/>
- <http://www.openstack.org/>
- <https://wiki.openstack.org/wiki/Rally>
- <https://wiki.openstack.org/wiki/Docker>
- <http://devstack.org/>
- [http://www.linux-kvm.org/page/Main\\_Page](http://www.linux-kvm.org/page/Main_Page)
- <https://github.com/stackforge/nova-docker>
- <https://github.com/dotcloud/docker-registry>
- <http://www.netperf.org/netperf/>
- <http://www.tokutek.com/products/iibench/>
- <http://www.brendangregg.com/activebenchmarking.html>
- <http://wiki.openvz.org/Performance>