

Egypten

Osaka

Eivind Uggedal

Hype

"Best features of Docker:

- Easy to get a project to the top of HackerNews
- New Github stars for old ideas
- Has lots of Twitter followers"

twitter.com/hipsterhacker/status/577587054642536449

Containere

- o Namespace i en kernel for at ulike samlinger av prosesser kan kjøre og kontrolleres uavhengig av hverandre:
 - o prosess (PID)
 - o nettverk
 - o mount
 - o brukere (UID/GID)
 - o ...

Tidslinje

- o 1967: VM - IBM CP-40
- o 1979: chroot(2) - Version 7 AT&T UNIX
- o 2000: jail(2) - FreeBSD 4.0
- o 2001: VServer - Linux patchset
- o 2005: zones - Solaris 10
- o 2005: OpenVZ - Linux patchset
- o 2007: WPARs - IBM AIX 6.1
- o 2007: SRP - HP-UX 11i
- o 2008: LXC - Linux
- o 2009: Heroku - PaaS
- o 2010: dotCloud - PaaS
- o 2013: Docker - Linux

Linux

- o Utvidelser i clone(2) flag
 - o 2002: CLONE_NEWNS - mount
 - o 2006: CLONE_NEWUTS - hostname
 - o 2006: CLONE_NEWIPC - SysV IPC / POSIX MQ
 - o 2008: CLONE_NEWPID - prosess
 - o 2009: CLONE_NEWNET - nettverk
 - o 2012: CLONE_NEWUSER - bruker
- o Control Groups (cgroups)
 - o 2008: cpuset - pinning av CPU/minne-enheter
 - o 2008: memory - begrense minne
 - o 2008: devices - hviteliste enheter
 - o 2010: blkio - disk I/O begrensnig
 - o 2012: net_prio - prioritering av nettverstrafikk
 - o 2015: pids - antall prosesser

Docker arkitektur

- o Master prosess, kjørt som root
- o Kontrolleres via RPC HTTP grensesnitt
 - o Unix domain socket
 - o TCP
- o Klient (fri flyt ved tilgang til socket/TCP)

PID 1

- o Upstream anbefaling
 - o Din app som PID 1
 - o Ikke supervisord/runit/s6/etc
- o Hva med zombie prosesser?
 - o Forby child prosessser?
- o Trenger du et OS
 - o Bruk LXC og la init(8) bli PID 1
 - o Få med cron(8) og syslogd(8) på kjøpet

Distroer

- o Mainstream distroer ikke designet for containere:
 - o Node: 642.2 MB
 - o Java: 641.9 MB
 - o Ubuntu: 187.9 MB
 - o CentOS: 172.3 MB
 - o Debian: 125.1 MB
- o Alternativ:
 - o Alpine: 5.2 MB
 - o Bygg din egen
 - o buildroot.net
 - o custom

Sikkerhet

- o Ressursoptimalisering - ikke sikkerhetsisolasjon!
 - o Linux kernel
 - o 4.3: 20 621 558 linjer
 - o Docker root daemon
- o Sikkerhetsoppdateringer

User namespace

- o Fjerner nødvendighet av root
- o Sikrere
- o Enklere
- o LXC: støttet i over 2år
- o Docker: eksperimentell (ubrukelig) versjon tilgjengelig forrige måned

Tillit

o hub.docker.com...

WTF

- o --privileged=true

Ytelse

- o Kjappere enn VM?
 - o clearlinux.org - VM booter på under 150ms
 - o Hardware virtualisering + virtio er generelt raskere på IO/net mens litt tregere på ren CPU ytelse
- o Kjappere enn LXC?
 - o Docker trenger ikke starte init med oppstartsscript
 - o I bredre distroer er disse optimalisert for container bruk

Effektivitet

- o Lagdeling og underliggende FS tenknologi (COW)
tregere og bruker ofte mer plass enn naiv en/to-lags
implementasjon (LXC)
- o Manglende garbage collection

Community

- o Bidro selv i starten
- o 1000+ åpne issues
- o \$\$\$
 - o RedHat
 - o Microsoft
- o CoreOS drama
 - o rkt
 - o opencontainers.org

"All problems in computer science can be solved by
another level of indirection...
Except for the problem of too many layers of indirection."

David Wheeler

Misbruk

"Containers will not fix your broken architecture.
You are welcome."

twitter.com/littleidea/status/659445920954642432

"Not even two dockers could contain the shittiness
of this application."

twitter.com/sadserver/status/625663479370883073

Egnethet

- o Ikke putt alt i en docker container
- o Hvorfor må din app contains?
 - o Er den for kompleks til å håndteres uten?
- o Dockerizing
 - o Antipattern
- o Kveg vs kjæledyr
 - o Er din app klar til å bli skutt i hodet på tilfeldige tidspunkt?
- o Stateless vs stateful
- o Distribuert system vs failover
- o Docker gjør ikke tar(1) obselete

- o Container specifications
- o Container runtimes
- o Container management
- o Container definition
- o Registries
- o Container OS
- o VM management
- o Scheduling
- o Cluster definition
- o Security
- o Log management

- o Service discovery
- o Dynamic config management
- o Container orchestration platforms
- o Hosted container platforms
- o Container platform management
- o Container-based PaaS
- o Networking
- o Monitoring
- o Data layer
- o Continuous integration/deployment
- o Getting started aides

Open Container Spec, Open Container Initiative, Docker, CoreOS, AppC, runc, libcontainer, rkt, Docker Engine, Docker daemon, Docker client, rkt CLI, Docker Image, Dockerfile, ACI (App Container Image), Docker Registry, Amazon EC2 Container Registry, Docker Hub, Google Container Registry, Quay.io, Docker Trusted Registry, CoreOS Enterprise Registry, boot2docker, RancherOS, Project Atomic, Ubuntu Core "Snappy", SmartOS, Photon OS, Docker Machine, Hashicorp Vagrant, Hashicorp Otto, Docker Swarm, fleet, Chronos, Docker Compose, fleet unit file, etcd, Marathon, Hashicorp Consul, Apache ZooKeeper, confd, Consul Template, Apache Mesos, Kubernetes, Hashicorp Nomad, Amazon EC2 Container Service (ECS), Google Container Engine, Tutum, RedHat Openshift, Joyent Triton, Giant Swarm, ProfitBricks, Modulus, Project Orca, Rancer, ContainerShip, panamax, Shipyard, Joyent SmartDataCenter, Mesosphere DCOS, CoreOS Tectonic, Nirmata, ContainerShip Enterprise, StackEngine, AppFormix, Deis, Flynn, RedHat Openshift Origin, Cisco Mantl, Deis Dokku, Docker port expose, Docker linking, libnetwork, flannel, Weave, Calico, Docker stats API, sysdig, cAdvisor, Weave Scope, Sysdig Cloud, ClusterHQ Flocker, Docker logs, logspout, Shippable, Wercker, Twistlock, Scarlock, Conjur, Docker Kinematic

Kompleksitet

o Docker engine: 580 511 linjer

Alternativer

- o LXC
- o Lettvektsløsninger:
 - o unshare(1)
 - o github.com/arachsys/containers
- o Direkte i app server (uWSGI)

```
#include ...

int main(int argc, char argv)
{
    uid_t uid = getuid();
    uid_t gid = getgid();

    unshare(CLONE_NEWUSER | CLONE_NEWNS | CLONE_NEWPID | CLONE_NEWNET);

    int fd = open("/proc/self/uid_map", O_RDWR);
    dprintf(fd, "%u %u 1\n", 0, uid);
    close(fd);
    fd = open("/proc/self/gid_map", O_RDWR);
    dprintf(fd, "%u %u 1\n", 0, gid);
    close(fd);

    setgroups(0, 0);

    chdir(argv[1]);
    mount("/dev", "./dev", 0, MS_BIND | MS_REC, 0);
    mount("/proc", "./proc", 0, MS_BIND | MS_REC, 0);
    chroot(".");

    execv(argv[2], argv+2);
}
```

Unikernels

- o Fremtiden?
- o Fjerner deler av OS i motsetning til å legge til et lag over
 - o Mindre
 - o Sikrere
 - o Kjappere
 - o Mer effektivt
- o Rumprun
 - o NetBSD
- o MirageOS
 - o OCaml

"The most revered feature of the modern operating system is support for running existing applications. Minimally implemented application support is a few thousand lines of code plus the drivers, as we demonstrated with the Rumprun unikernel. Therefore, there is no reason to port and cram an operating system into every problem space."

Antti Kantee. "The Rise and Fall of the Operating System".
USENIX ;login: October 2015, Vol. 40, No. 5

Slides

- o <http://git.uggedal.com/presentations/tree/boycott-docker>