

LINUX CHECKPOINT/RESTART

FACCIAMOCI UN GIRO SULLA MACCHINA DEL TEMPO!

Micky Del Favero
micky@linux.it

BLUG - Belluno Linux User Group

LinuxDay 2011 - Belluno 22 ottobre 2011

DEFINIZIONE DI CHECKPOINT/RESTART

Checkpoint salvataggio uno stato arbitrario di un processo

DEFINIZIONE DI CHECKPOINT/RESTART

Checkpoint salvataggio uno stato arbitrario di un processo

Restart ripartenza di un processo da uno stato che un *checkpoint* ha precedentemente salvato

COME FUNZIONA?

Salva e ripristina uno stato di un processo in esecuzione

- ▶ file
- ▶ memoria
- ▶ IPC
- ▶ socket
- ▶ pty
- ▶ etc.

COME FUNZIONA?

Salva e ripristina uno stato di un processo in esecuzione

Opera a livello OS, granularità pari a task / container

COME FUNZIONA?

Salva e ripristina uno stato di un processo in esecuzione

Opera a livello OS, granularità pari a task / container

Simile a snapshot di VM, ma

- ▶ per processo in luogo di per sistema
- ▶ può essere più flessibile
- ▶ è più complicato
- ▶ non richiede supporto hardware e/o hardware virtualizzato

A COSA PUÒ SERVIRE?

Ripristino del processo ad uno stato precedente.

- ▶ Tolleranza ai guasti
 - ▶ Elaborazioni di lunga durata
 - ▶ HPC / Cloud
- ▶ Debug / Collaudo
 - ▶ ripartenza da stato stabile
 - ▶ ripetizione di passi specifici
 - ▶ distribuzione su più host
- ▶ Lancio rapido di applicazioni
- ▶ Macchina del tempo :-)

A COSA PUÒ SERVIRE?

Ripristino del processo ad uno stato precedente.

Sospensione di un processo e successiva ripartenza, anche su host diversi.

- ▶ Gestione della memoria
- ▶ Gestione del sistema
 - ▶ gestione tempo macchina
 - ▶ sospensione di applicazioni secondarie in caso di carico eccessivo
- ▶ Gestione sessione utente
 - ▶ sessione mobile
 - ▶ VPS / container

A COSA PUÒ SERVIRE?

Ripristino del processo ad uno stato precedente.

Sospensione di un processo e successiva ripartenza, anche su host diversi.

Migrazione di processi su host remoti.

- ▶ Bilanciamento del carico
 - ▶ HPC
 - ▶ Cloud
 - ▶ VPS
- ▶ Zero-downtime grazie a migrazione dei processi su server remoti
- ▶ Alta affidabilità
 - ▶ Backup live
 - ▶ Ripristino in caso di crash

IMPLEMENTAZIONI

Distributed MultiThreaded CheckPointing (DMTCP)

- ▶ Checkpoint/Restart completamente in userspace
- ▶ Basato su *LD_PRELOAD*
- ▶ Non richiede nessuna patch!

IMPLEMENTAZIONI

Distributed MultiThreaded CheckPointing (DMTCP)

Qemu / KVM

- ▶ Checkpoint/Restart in userspace
- ▶ Granularita pari ad un intero OS

IMPLEMENTAZIONI

Distributed MultiThreaded CheckPointing (DMTCP)

Qemu / KVM

Berkeley Lab C/R (BLCR)

- ▶ Approccio ibrido kernel/userspace
- ▶ Richiede patch del kernel
- ▶ Progetto abbandonato?

IMPLEMENTAZIONI

Distributed MultiThreaded CheckPointing (DMTCP)

Qemu / KVM

Berkeley Lab C/R (BLCR)

OpenVZ

- ▶ Checkpoint/Restart in kernel space
- ▶ Richiede patch del kernel

IMPLEMENTAZIONI

Distributed MultiThreaded CheckPointing (DMTCP)

Qemu / KVM

Berkeley Lab C/R (BLCR)

OpenVZ

Linux-CR

- ▶ Checkpoint/Restart in kernel space con interfaccia di Restart in userspace
- ▶ Branch del kernel (*git clone...*)
- ▶ Usa cgroups

IMPLEMENTAZIONI

Distributed MultiThreaded CheckPointing (DMTCP)

Qemu / KVM

Berkeley Lab C/R (BLCR)

OpenVZ

Linux-CR

DMTCP

DMCTP DISTRIBUTED MULTITHREADED CHECKPOINTING

Completamente in userspace

Non richiede privilegi particolari

Non richiede moduli o modifiche al kernel né al software da controllare

Usato nell'esperimento CMS del Large Hadron Collider del CERN

STATO DELL'ARTE

Gestisce automaticamente

- ▶ processi: fork, exec, ssh, pid, etc
- ▶ ipc: mutex, sem
- ▶ socket: unix, ipv4, ipv6
- ▶ terminali: pipe, pty, terminal mode, signal handler
- ▶ file: I/O, ownership, file condivisi, shm (via mmap)

STATO DELL'ARTE

Supporta una vasta gamma di applicazioni, anche binarie, fra cui

- ▶ Matlab
- ▶ Perl, Python, shell
- ▶ VNC
- ▶ OpenMPI

STATO DELL'ARTE

Sistemi non supportati

- ▶ Infiniband
- ▶ Myrinet

FUNZIONAMENTO

Lancio del *coordinator*

Lancio del processo in *dmtcp*

Checkpoint (Salvataggio dell'immagine)

Termine del processo

Eventuale *Restart* da un *Checkpoint* precedente.

ESEMPIO PRATICO

Lancio del coordinator

ESEMPIO PRATICO

```
$ dmtcp_coordinator
```

ESEMPIO PRATICO

```
$ dmtcp_coordinator
```

Lancio del processo in dmtcp

ESEMPIO PRATICO

```
$ dmtcp_coordinator
```

```
$ dmtcp_checkpoint command -par0 ... -parN
```

ESEMPIO PRATICO

```
$ dmtcp_coordinator
```

```
$ dmtcp_checkpoint command -par0 ... -parN  
Checkpoint
```

ESEMPIO PRATICO

```
$ dmtcp_coordinator
```

```
$ dmtcp_checkpoint command -par0 ... -parN
```

```
> c
```

ESEMPIO PRATICO

```
$ dmtcp_coordinator
```

```
$ dmtcp_checkpoint command -par0 ... -parN
```

```
> c
```

Il processo termina

ESEMPIO PRATICO

```
$ dmtcp_coordinator
```

```
$ dmtcp_checkpoint command -par0 ... -parN
```

```
> c
```

Il processo termina

Restart

ESEMPIO PRATICO

```
$ dmtcp_coordinator
```

```
$ dmtcp_checkpoint command -par0 ... -parN
```

```
> c
```

Il processo termina

```
$ ls
```

```
ckpt_command_1fa17adc8b68249.dmtcp
```

```
dmtcp_restart_script_1fa17adc8b68249.sh
```

```
dmtcp_restart_script.sh
```

```
$ ./dmtcp_restart_script.sh
```

Linux-CR

LINUX CR

Progetto iniziato da Oren Laadan e IBM

Punta ad essere integrato nel kernel ufficiale

Usato nei supercomputer PERCS / BlueWater

OBIETTIVI

Trasparenza da e verso i processi

OBIETTIVI

Trasparenza da e verso i processi

Affidabilità

- ▶ Checkpoint ok \Rightarrow Restart
- ▶ Segnala situazioni non *Checkpoint*-abili
- ▶ Checkpoint non invasivo

OBIETTIVI

Trasparenza da e verso i processi

Affidabilità

Performance

- ▶ non ha impatto sulle performance del sistema
- ▶ il codice ha un'immagine regionevoles

OBIETTIVI

Trasparenza da e verso i processi

Affidabilità

Performance

Mantenibilità

- ▶ codice generico in kernel/checkpoint
- ▶ codice dei sottosistemi in subsystem/checkpoint
- ▶ cambio versione non provoca grossi impatti (2.6.27 in avanti)

STATO DELL'ARTE

Sottosistemi supportati

- ▶ tasks: thread, segnali, credenziali, etc.
- ▶ tutti i namespace tranne mount-ns
- ▶ sysvipc: shm, msg, sem
- ▶ file: regolari, fifo/pipe, epoll, eventi, dispositivi
- ▶ socket: unix, ipv4, ipv6
- ▶ sicurezza: smack, selinux

STATO DELL'ARTE

Sottosistemi supportati

Sottosistemi non supportati

- ▶ lock, leases, owner
- ▶ file cancellati
- ▶ inotify/dnotify
- ▶ mount
- ▶ filesystem /proc

STATO DELL'ARTE

Sottosistemi supportati

Sottosistemi non supportati

Architetture supportate

- ▶ i386
- ▶ x86-64
- ▶ s390x
- ▶ PowerPC
- ▶ ARM

FUNZIONAMENTO: CHECKPOINT

1. Congela la gerarchia dei processi
2. Salva i dati
3. Salva la gerarchia dei processi
4. Salva lo stato di tutti i task
5. Scongela (o killa :-) la gerarchia dei processi

FUNZIONAMENTO: RESTART

1. Crea il container
2. Crea la gerarchia dei processi
3. Ripristina lo stato di tutti i task
4. Riprendi l'esecuzione

INSTALLAZIONE

```
git clone git://www.linux-cr.org/pub/git/linux-cr
git clone git://www.linux-cr.org/pub/git/user-cr
git clone git://www.linux-cr.org/pub/git/tests-cr
cp /boot/config-`uname -r` .config
$EDITOR. .config
make-kpkg --intrd kernel_image
cd user-cr && make all
cd test-cr && make all
```

ESEMPIO PRATICO

```
mkdir /cgroup/mytask  
echo $PID > /cgroup/mytask/tasks  
echo FROZEN > /cgroup/mytask/freezer.state  
checkpoint $PID > statefile  
echo THAWED > /cgroup/mytask/freezer.state  
restart --pidns < statefile
```

DOCUMENTAZIONE, SORGENTI, ETC.

DMTCP

<http://dmtcp.sourceforge.net/>

Linux-CR

<http://www.linux-cr.org/>

<http://www.linux-cr.org/pub/git/linux-cr>

<http://www.linux-cr.org/pub/git/user-cr>

<http://www.linux-cr.org/pub/git/tests-cr>

Domande?

Grazie per l'attenzione!



Attribuzione Non commerciale Condividi allo stesso modo 3.0 Italia CC BY-NC-SA 3.0)