

DE LA RECHERCHE À L'INDUSTRIE



www.cea.fr

Debugging kernel

Outils pour comprendre ce qu'il se passe la dedans

Dominique Martinet

CEA

January 9, 2020

- 1 Intro
 - Préface
 - Setup du TP
- 2 Crash
- 3 Perf
- 4 SystemTap
- 5 eBPF : bcc-tools, bpftrace

Intro



Crash



Perf



SystemTap



eBPF : bcc-tools, bpftrace



- Formation recettes de cuisines, principalement sur crash
- Pas une solution magique, il faut toujours :
 - Chercher à comprendre ce qu'on analyse
 - Lire du code
 - Jouer à Sherlock Holmes
- Mais vous devriez ressortir d'ici avec au moins un outil de plus pour tout ça.

Intro



Crash



Perf



SystemTap



eBPF : bcc-tools, bpftrace



■ Bases de C

- Bases d'assembly peut aider, mais pas vu ici

■ Commandes unix standard

■ Bon à connaitre (revus rapidement lors du TP) :

- Savoir ce qu'est une stack
- Connaitre un peu les types de base du kernel (en particulier listes)

- Une image de VM est disponible sur simple mail, n'hésitez pas à demander par email. Elle est nécessaire pour reproduire les TPs.
- Importez l'image disque dans un outil tout fait ou lancez une commande semblable à la suivante :

```
$ qemu-system-x86_64 -hda kdebug.qcow2 -serial mon:stdio \
    -enable-kvm -nographic -smp 16 -m 20G \
    -cdrom ~/cloud.iso
```

- En l'absence d'image cloud-init, vous devriez avoir un prompt de login où vous pourrez vous connecter avec root :toor.



Intro



Crash



Perf



SystemTap



eBPF : bcc-tools, bpftrace



- Se connecter sur les postes si ce n'est pas déjà fait.
- Se connecter à inti avec les mêmes identifiants
- Pour pouvoir copier/coller des choses si besoin ou revenir sur une ancienne slide, celles-ci sont disponible sur `~stag01/presentation.fr.pdf`

Démarrer une VM.

Le script génère une clé ssh, crée un script cloud-init avec cette clé et lance la VM pcocc préparée pour la formation. La VM est un simple linux ocean avec les packets debuginfo installés, les outils utilisés, et des crash dumps copiés d'ailleurs.

```
$ module load datadir/formation
$ $FORMATION_CCCSCRATCHDIR/start_debug_vm.sh
```

Usage

- pcocc ssh root@vm0
- pcocc scp root@vm0:foo bar
- pcocc console vm0, pcocc reset vm0, etc etc...

Liste des paquets

lustre lustre-debuginfo kernel-debuginfo crash systemtap perf bcc-tools gdb cscope
ripgrep bash-completion htop bpftool (llvm-toolset-7.0 devtoolset-8 bpfftrace)

1 Intro

2 Crash

- LBUG client lustre
- Charge serveur lustre
- Coup d'œil sur un autre exemple

3 Perf

4 SystemTap

5 eBPF : bcc-tools, bpftrace

- Outil pour analyser l'état du kernel
- Principalement un gros wrapper autour de gdb

Live

- Ne peut pas donner les stacks des threads actifs (comme /proc/pid/stack)
- Selon le driver peut modifier la mémoire de la machine (pas par défaut)
- Ne prend pas de locks donc attention aux choses qui bougent

Post-mortem (kexec-kdump)

- dans /var/crash/<ip>-<date>/vmcore

Lancer crash

```
# cd /var/crash/client_lbug
# crash /usr/lib/debug/lib/modules/3.10.0*/vmlinuz vmcore
...
PANIC: "Kernel panic - not syncing: LBUG"
PID: 12295
COMMAND: "trinity-c30"
...
crash>
```

Notes crash

- Les infos de la machine données ici peuvent être réaffichées si besoin avec `sys`
- Le process listé est celui qui était actif sur le cpu qui a causé le crash
- Ce process est utilisé pour toutes les commandes qui prennent un pid/task si on n'en précise pas
- En cas d'urgence, `help / help <command>` sont de bons amis

Intro Crash Perf SystemTap eBPF : bcc-tools, bpfftrace

○○○○○○○ ○○○○○●○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○ ○○○○○○○○○○ ○○○○○○○○ ○○○○○○○○○○○○

```
crash> bt -FF
...
#6 [ffff9f94f45dfe80] do_vfs_ioctl at ffffffff94e569d0
ffff9f94f45dfe88: [ffff9f94e6411cc0:dentry] [ffff9f94f45d669f:mm_struct]
ffff9f94f45dfe98: ffff9f9f958ed1590b ffff9f94f45dfe80
ffff9f94f45dfea8: do_setitimer+476 0000000000000000
ffff9f94f45dfef8: 00000000ecef1f068 [ffff9f956da1ff00:kmalloc-256]
ffff9f94f45dfef8: 0000000000000037 00000000ecef1f068
ffff9f94f45dfef8: [ffff9f956da1ff00:kmalloc-256] 0000000000000037
ffff9f94f45dfef8: 000000004008669f 000000007fffff00
ffff9f94f45dfef8: ffff9f94f45dff48 sys_ioctl+161
```

Notes crash

- On retrouve la dentry et le file de tout à l'heure
- la syntaxe [addr:slab] n'indique que l'appartenance à un pool d'allocation : ce n'est pas nécessairement le début de la structure !
- les pools de même taille sont regroupés, e.g. même si le code a un kmem_cache_create() on n'aura pas forcément l'indice

```

crash> mod -s libcfs
crash> mod -s lustre
crash> bt -l

...
#4 [ffff9f94f45dfd50] ll_put_grouplock at ffffffff0fcb92c [lustre]
  /usr/src/debug/lustre-2.12.3/lustre/llite/file.c: 2200
#5 [ffff9f94f45dfda0] ll_file_ioctl at ffffffff0fde8d6 [lustre]
  /usr/src/debug/lustre-2.12.3/lustre/llite/file.c: 3388
#6 [ffff9f94f45dfe80] do_vfs_ioctl at ffffffff4a5e69d0
  /usr/src/debug/kernel-3.10.0-957.21.3.el7/linux-3.10.0-957.21.3.el7.x86_64/fs/ioctl.c:
...

```

Notes crash

- mod -s <module> **ou** mod -S
- ... reste à lire le code et comprendre ce qu'il s'est passé !
- cscope est votre ami

Intro



Crash



Perf



SystemTap



eBPF : bcc-tools, bpftrace



- Crash n'est pas un outil magique, il faut vraiment regarder dans le code ce qui peut nous faire arriver là.

- Crash n'est pas un outil magique, il faut vraiment regarder dans le code ce qui peut nous faire arriver là.

```
@@ static int ll_put_grouplock(struct inode *inode, ...
    if (!(fd->fd_flags & LL_FILE_GROUP_LOCKED)) {
        spin_unlock(&lli->lli_lock);
        CWARN("no group lock held\n");
        RETURN(-EINVAL);
    }

    LASSERT(fd->fd_grouplock.lg_lock != NULL);
```

- Crash n'est pas un outil magique, il faut vraiment regarder dans le code ce qui peut nous faire arriver là.

```
@@ static int ll_put_grouplock(struct inode *inode, ...
    if (!(fd->fd_flags & LL_FILE_GROUP_LOCKED)) {
        spin_unlock(&lli->lli_lock);
        CWARN("no group lock held\n");
        RETURN(-EINVAL);
    }

    LASSERT(fd->fd_grouplock.lg_lock != NULL);
```

- Quelque chose a remis `fd_grouplock` à zero sans changer le flag ?
- Ou bien `fd_flags` a été modifié sans toucher au grouplock ?

- Crash n'est pas un outil magique, il faut vraiment regarder dans le code ce qui peut nous faire arriver là.

```
@@ static int ll_put_grouplock(struct inode *inode, ...
    if (!(fd->fd_flags & LL_FILE_GROUP_LOCKED)) {
        spin_unlock(&lli->lli_lock);
        CWARN("no group lock held\n");
        RETURN(-EINVAL);
    }

    LASSERT(fd->fd_grouplock.lg_lock != NULL);
```

- Quelque chose a remis `fd_grouplock` à zero sans changer le flag ?
- Ou bien `fd_flags` a été modifié sans toucher au grouplock ?
- ... Il y a autre ioctl qui permet de mettre des valeurs arbitraires dans `fd->fd_flags` sans aucune vérification
- LU-13257 ouvert pour l'occasion (Bug trouvé pour cette formation !)


```
$ cd /var/crash/server_load
$ crash /usr/lib/debug/lib/modules/3.10.0*/vmlinux vmcore
...
    PANIC: "SysRq : Trigger a crash"
        PID: 4715
        COMMAND: "bash"
...
crash>
```

- ### Notes crash
- Pas de crash franc : SysRq crash
 - Il faut faire un peu plus le tour de la machine

```
crash> ps
  PID   PPID  CPU   TASK                               ST  %MEM  VSZ   RSS  COMM
  > 0     0     0     ffffffff9da18480                  RU  0.0   0     0   [swapper/0]
  > 0     0     1     ffff9177f60a4100                  RU  0.0   0     0   [swapper/1]
  ...
  5387   2     4     ffff917cbba30000                  IN  0.0   0     0   [ll_ost02_006]
  5388   2     5     ffff917cbba36180                  IN  0.0   0     0   [mdt02_003]
  > 5389   2     0     ffff917cbba330c0                  RU  0.0   0     0   [ll_ost00_003]
  ...
  5532   2     3     ffff917c6a9b1040                  RU  0.0   0     0   [ll_ost_io01_006]
  5533   2     6     ffff917c6a9b2080                  IN  0.0   0     0   [mdt03_011]
  5534   2     0     ffff917c6a9b30c0                  IN  0.0   0     0   [mdt00_015]
  5535   2     7     ffff917c6a9b4100                  IN  0.0   0     0   [mdt03_012]
  5536   2     5     ffff917c6a9b5140                  UN  0.0   0     0   [ll_ost_io02_006]
  ...
```

Notes crash

>, IN, UN, RU Active, Interruptible, Uninterruptible, Runnable
 ZO, ST, TR, DE Zombie, Stopped, Tracing, Dead
 SW, WA, PA, NE, EX... Swapping, Waking, Parked, New, Exclusive

```
crash> ps -l
[1834107714279] [IN] PID: 4831 TASK: ffff9..b230c0 CPU: 7 COMMAND: "socknal_sd03_01"
[1834107636683] [IN] PID: 4828 TASK: ffff9..b25140 CPU: 5 COMMAND: "socknal_sd02_00"
[1834107619697] [RU] PID: 6121 TASK: ffff9..8f6180 CPU: 2 COMMAND: "ll_ost_io01_042"
[1834107601107] [IN] PID: 4829 TASK: ffff9..b22080 CPU: 4 COMMAND: "socknal_sd02_01"
...
[1833607166340] [RU] PID: 6032 TASK: ffff9..f65140 CPU: 3 COMMAND: "ll_ost_io01_023"
[1833600488498] [UN] PID: 6225 TASK: ffff9..b44100 CPU: 2 COMMAND: "ll_ost_io01_068"
[1833568488187] [RU] PID: 5532 TASK: ffff9..9b1040 CPU: 3 COMMAND: "ll_ost_io01_006"
[1833531745921] [RU] PID: 6129 TASK: ffff9..208000 CPU: 3 COMMAND: "ll_ost_io01_044"
[1833522652478] [RU] PID: 6239 TASK: ffff9..aec100 CPU: 3 COMMAND: "ll_ost_io01_074"
[1833501019679] [RU] PID: 6094 TASK: ffff9..1a0000 CPU: 3 COMMAND: "ll_ost_io01_035"
[1833456118992] [RU] PID: 6214 TASK: ffff9..2b8000 CPU: 3 COMMAND: "ll_ost_io01_066"
```

Notes crash

- première colonne entre crochet : dernière date à laquelle le thread a été schedulé (uptime en nanoseconde)
- processus triés sur cette colonne : utile pour trouver des groupes de threads bloqués
- 1834107714279 – 1833456118992 = 652ms
 - loin des traces à partir de 120s en UN ou RU sans tourner
- dmesg | tail -> [1834.107266] ...

```
crash> foreach bt > foreach_bt
crash> ! awk '/PID:/ {PID=$2;PNAME=$8}
           /^ ?#/ {print PNAME": "$1" "$3" "$5}' foreach_bt |
           clubak -c

-----
"ll_ost00_003"
-----
...
#7 queued_spin_lock_slowpath ffffffff9d55d28b
#8 _raw_spin_lock ffffffff9d56b760
#9 lock_res_and_lock ffffffff0fcc02c
#10 ldlm_work_bl_ast_lock ffffffff0fcf0d9
#11 ptlrpc_check_set ffffffff01014030
#12 ptlrpc_check_set ffffffff0101587b
#13 ptlrpc_set_wait ffffffff01015c24
#14 ldlm_run_ast_work ffffffff0fd3055
#15 ldlm_handle_conflict_lock ffffffff0fd3780
#16 ldlm_lock_enqueue ffffffff0fd3cc3
#17 ldlm_handle_enqueue0 ffffffff0ffc336
```

Notes crash

- plein d'arguments possibles : `regex ('mdt0.*')`, `state (active, RU, UN)`...
- `foreach bt -R lock_res_and_lock, foreach bt -FF, foreach files...`

Intro Crash Perf SystemTap eBPF : bcc-tools, bpfftrace
 ○○○○○○ ○○○○○○○○○○○○○○○○○●○○○○○○○○○○○○○○○○○○○ ○○○○○○○○○○ ○○○○○○○○ ○○○○○○○○○○○○○○

```

crash> mod -S > /dev/null
crash> bt -FF 5389

PID: 5389  TASK: ffff917cbba330c0  CPU: 0  COMMAND: "ll_ost00_003"
...
#14 [ffff917c9eb33b40] ldlm_run_ast_work at ffffffff0fd3055 [ptlrpc]
ffff917c9eb33b48: [ffff917c62522f40:kioctx] [ffff917c981c3b40:kmallo-192]
ffff917c9eb33b58: [ffff917c981c3b5c:kmallo-192] ffff917c9eb33c68
ffff917c9eb33b68: ffff917c9eb33ba0 ldlm_handle_conflict_lock+112
#15 [ffff917c9eb33b70] ldlm_handle_conflict_lock at ffffffff0fd3780 [ptlrpc]
ffff917c9eb33b78: ffff917c9eb33c68 0000000000000000
ffff917c9eb33b88: ffff917c9eb33bd8 [ffff917c981c3b40:kmallo-192]
ffff917c9eb33b98: [ffff917c62522f40:kioctx] ffff917c9eb33c18
ffff917c9eb33ba8: ldlm_lock_enqueue+723
#16 [ffff917c9eb33ba8] ldlm_lock_enqueue at ffffffff0fd3cc3 [ptlrpc]
ffff917c9eb33bb0: 000000008058614a [ffff917cb168a430:kmallo-2048]
ffff917c9eb33bc0: ffff91772a3780e0 ldlm_process_extent_lock
ffff917c9eb33bd0: 00000000c11e5a0 [ffff917c59d6be88:kioctx]
ffff917c9eb33be0: [ffff917c6059e808:kioctx] 000000008058614a
ffff917c9eb33bf0: [ffff917cb168a050:kmallo-2048] ffff91772a3780e0
ffff917c9eb33c00: 0000000000000000 [ffff917cdd6a2000:kmallo-1024]
ffff917c9eb33c10: [ffff917cf510fc40:kmallo-64] ffff917c9eb33ca8
ffff917c9eb33c20: ldlm_handle_enqueue0+2646
#17 [ffff917c9eb33c20] ldlm_handle_enqueue0 at ffffffff0ffc336 [ptlrpc]
ffff917c9eb33c28: ffff917700000000 ffff917700000000
ffff917c9eb33c38: lustre_swab_ldlm_request [ffff917cb168a430:kmallo-2048]
ffff917c9eb33c48: 0000006800000001 [ffff917c62522f40:kioctx]
ffff917c9eb33c58: 0000000000000038 [ffff917cb168a430:kmallo-2048]
ffff917c9eb33c68: 0000000000020000 [ffff917c62522f40:kioctx]
ffff917c9eb33c78: 000000008058614a [ffff917c5ead9960:sigqueue]
ffff917c9eb33c88: [ffff917cb168a050:kmallo-2048] [ffff917ca0e6a110:kmallo-2048]
ffff917c9eb33c98: [ffff917cb168a050:kmallo-2048] tgt_dlm_handlers
  
```

Intro
○○○○○Crash
○○○○○○○○○○○○○○○○●○○○○○○○○○○○○○○○○○○○○Perf
○○○○○○○○○○○○SystemTap
○○○○○○○○○○eBPF : bcc-tools, bpftrace
○○○○○○○○○○○○○○○○

```

crash> set radix 16
crash> struct -o ldlm_resource

struct ldlm_resource {
  [0x0] struct ldlm_ns_bucket *lr_ns_bucket;
  [0x8] struct hlist_node lr_hash;
  [0x18] atomic_t lr_refcount;
  [0x1c] spinlock_t lr_lock;
  [0x20] struct list_head lr_granted;
  [0x30] struct list_head lr_waiting;
  [0x40] struct ldlm_res_id lr_name;
  ...
  [0x70] enum ldlm_type lr_type;
  [0x74] int lr_lvb_len;
  [0x78] struct mutex lr_lvb_mutex;
  [0xa0] void *lr_lvb_data;
  [0xa8] _Bool lr_lvb_initialized;
  [0xa9] struct lu_ref lr_reference;
}
SIZE: 0xb0

```

Notes crash

- 0xb0 = 176 : chercher kmalloc-192
- Une seule valeur qui tombe rond dans bt -FF la dedans...

```
crash> struct ldlm_resource ffff917c981c3b40

  lr_ns_bucket = 0xffffb40a6eald018,
  ...
  lr_refcount = {
    counter = 0x145d7
  },
  ...
  lr_granted = {
    next = 0xffff917c6648f420,
    prev = 0xffff91772946ble0
  },
  lr_waiting = {
    next = 0xffff917c62522fa0,
    prev = 0xffff917c867b0060
  },
  lr_name = {
    name = {0x28, 0x0, 0x0, 0x0}
  },
  ...
  lr_type = LDLM_EXTENT,
  lr_lvb_len = 0x38,
```

Notes crash

- `lr_name` souvent un fid, mais ici sur un OST on a l'oid
- `lr_refcount` Nombre de références à la ressource
- `lr_granted/waiting` listes de locks

Intro
○○○○○Crash
○○○○○○○○○○○○○○○○○○●○○○○○○○○○○○○○○○○○○Perf
○○○○○○○○○○○○SystemTap
○○○○○○○○○○eBPF : bcc-tools, bpftrace
○○○○○○○○○○○○○○○○○○

```

crash> struct ldlm_resource | grep lr_ns_bucket
    struct ldlm_ns_bucket *lr_ns_bucket;
crash> struct ldlm_ns_bucket 0xffffb40a6ea1d018
    nsb_namespace = 0xffff917cdd6a2000,
    ...
crash> struct ldlm_ns_bucket | grep nsb_namespace
    struct ldlm_namespace *nsb_namespace;
crash> struct ldlm_namespace 0xffff917cdd6a2000
    ...
    ns_name = 0xffff917ce5b5ae60 "filter-testfs0-OST0000_UUID",

```

Notes crash

- Déréférences consécutives non triviales mais possibles

Intro

○○○○○

Crash

○○○○○○○○○○○○○○○○○○○○●○○○○○○○○○○○○○○○○○○

Perf

○○○○○○○○○○○○

SystemTap

○○○○○○○○○○

eBPF : bcc-tools, bpftrace

○○○○○○○○○○○○○○○○○○

```

# mount -o loop -t ldiskfs /tmp/ost0 /media/
# ls /media/O/0/d$( (40%32) )/40
/media/O/0/d8/40
# umount /media
# debugfs -c /tmp/ost0
...
debugfs: stat /O/0/d8/40
...
    fid: parent=[0x200000401:0x4d:0x0] stripe=0
        stripe_size=1048576 stripe_count=1
...
^D
# lfs fid2path /mnt/lustre 0x200000401:0x4d:0x0
/mnt/lustre/test

```

Intro
○○○○○Crash
○○○○○○○○○○○○○○○○○○○○●○○○○○○○○○○○○○○○○○○Perf
○○○○○○○○○○○○SystemTap
○○○○○○○○○○eBPF : bcc-tools, bpftrace
○○○○○○○○○○○○○○○○○○

```

crash> help list
crash> struct list_head
struct list_head {
    struct list_head *next;
    struct list_head *prev;
}
SIZE: 16

```

Notes

- double-linked list
- Souvent tête de liste séparée : `list -H`
 - i.e. `lustre/ldlm/ldlm_internal.h:38: extern struct list_head ldlm_srv_namespace_list;`
 - `/!\` même type pour la tête que le reste

```
crash> list -H ldlm_srv_namespace_list
ffff917cdd6a2000
...
```

Usage

```
list_for_each_entry(ns, ldlm_namespace_list(LDLM_NAMESPACE_SERVER),
                    ns_list_chain) {
```

```
crash> list -s ldlm_namespace ldlm_namespace.ns_list_chain
-H ldlm_srv_namespace_list
ffff917cdd6a2000
struct ldlm_namespace {
    ns_obd = 0xffff917ca0d620f0,
    ns_client = LDLM_NAMESPACE_SERVER,
    ns_name = 0xffff917ce5b5ae60 "filter-testfs0-OST0000_UUID",
    ns_rs_hash = 0xffff917ce0680000,
    ...
```

```

crash> struct -o ldlm_resource | grep -E 'lr_granted|lr_waiting'
[0x20] struct list_head lr_granted;
[0x30] struct list_head lr_waiting;
crash> struct -d ldlm_resource.lr_refcount ffff917c981c3b40
lr_refcount = {
    counter = 83415
}
crash> list -H ffff917c981c3b60 | wc -l
82718
crash> list -H ffff917c981c3b70 | wc -l
3

```

Notes crash

- ici (voir code) lr_granted/lr_waiting sont des têtes de listes à part
- Contiennent des struct ldlm_lock sur le champ l_res_link


```

crash> struct -o ldlm_resource

struct ldlm_lock {
...
[0x48] struct ldlm_resource *l_resource;
[0x60] struct list_head l_res_link;
[0x98] enum ldlm_mode l_req_mode;
[0x9c] enum ldlm_mode l_granted_mode;
[0xb8] struct obd_export *l_export;
[0x100] __u64 l_flags;
        union {
[0x160]     time64_t l_activity;
[0x160]     time64_t l_blast_sent;
        };
[0x1c0] __u32 l_pid;
[0x1f8] struct ldlm_lock *l_blocking_lock;
...
}
SIZE: 0x230

```

```
crash> list -H ffff917c981c3b70 ldlm_lock.l_res_link
          -s ldlm_lock.l_req_mode
ffff917c62522f40
  l_req_mode = LCK_PW
ffff917c641e5f80
  l_req_mode = LCK_PR
crash> list -H ffff917c981c3b60 ldlm_lock.l_res_link
          -S ldlm_lock.l_req_mode | grep -B 1 'l_req_mode = 0x2'
ffff917c6648f3c0
  l_req_mode = 0x2
ffff917c96fdfcc0
  l_req_mode = 0x2
...

```

- Notes crash**
- -s plus rapide mais pas de résolution des symboles
 - plusieurs LCK_PW en granted avec grandes majorités de PR, à priori incompatibles, mais l_req_extent différents.

Intro

○○○○○

Crash

○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○●○○○○○○○○○○

Perf

○○○○○○○○○○○○

SystemTap

○○○○○○○○○

eBPF : bcc-tools, bpftrace

○○○○○○○○○○○○○○

```
crash> list -H ffff917c981c3b70 ldlm_lock.l_res_link
          -S ldlm_lock.l_export
ffff917c62522f40
          l_export = 0xffff917c96f59000
```

```
crash> list -H ffff917c981c3b70 ldlm_lock.l_res_link
           -S ldlm_lock.l_export
ffff917c62522f40
    l_export = 0xffff917c96f59000
```

```
crash> struct obd_export.exp_connection 0xffff917c96f59000
    exp_connection = 0xffff917cdee0ba80
```

```
crash> list -H ffff917c981c3b70 ldlm_lock.l_res_link
           -S ldlm_lock.l_export
ffff917c62522f40
l_export = 0xffff917c96f59000
```

```
crash> struct obd_export.exp_connection 0xffff917c96f59000
exp_connection = 0xffff917cdee0ba80
```

```
crash> struct ptlrpc_connection 0xffff917cdee0ba80
struct ptlrpc_connection {
  c_peer = {
    nid = 0x200000ac80002,
    pid = 0x3039
  },
```

- ### Notes crash
- nid = lnd type, lnd number, ip/id
 - ici : tcp, 0, 0a.c8.00.02 = 10.200.0.2
 - o2ib : 0005001a0a800002
 - ptlf : 000f001500000014
 - net -N 0x200000ac80002 → 2.0.200.10

```

crash> list -H ffff917c981c3b60 ldlm_lock.1_res_link
-S ldlm_lock.1_export | awk '/l_export/ { print $3 }'
> export
crash> list -H ffff917c981c3b70 ldlm_lock.1_res_link
-S ldlm_lock.1_export | awk '/l_export/ { print $3 }'
>> export
crash> !head -n 1 export
0xffff917c96f59000
crash> struct obd_export.exp_connection 0xffff917c96f59000
exp_connection = 0xffff917cdee0ba80
crash> struct -o obd_export.exp_connection
struct obd_export {
  [0x118] struct ptlrpc_connection *exp_connection;
}
crash> rd -o 0x118 0xffff917c96f59000
ffff917c96f59118: ffff917cdee0ba80      ....|...
crash> rd -o 0x118 < export | awk '{ print $2 }' > conn
  
```

```

crash> struct -o ptlrpc_connection.c_peer
struct ptlrpc_connection {
  [0x18] struct lnet_process_id c_peer;
}
crash> !head conn
ffff917cdee0ba80
crash> struct ptlrpc_connection.c_peer.nid ffff917cdee0ba80
c_peer.nid = 0x200000ac80002,
crash> rd -o 0x18 ffff917cdee0ba80
ffff917cdee0ba98: 000200000ac80002 .....
crash> rd -o 0x18 < conn | awk '{ print $2 }' |
      sort | uniq -c | sort -h
      998 000200000ac8002c
      1019 000200000ac80041
      1024 000200000ac80024
      1029 000200000ac8002f
      ...
      1435 000200000ac80006
      1436 000200000ac8001e
      1445 000200000ac80016
  
```


Notes crash

- On utilise `rd` et pas `struct foo.bar` pour moins passer par `gdb` (beaucoup plus lent, un patch est en cours pour avoir quelque chose comme le `-S` vs `-s` de `list`)
- Abuser des fichiers et commandes externes
- L'extension `crash-extscript` pourrait aider si besoin, mais au final on s'y habitue...

Intro ○○○○○○ Crash ○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○○●○○○○○○ Perf ○○○○○○○○○○○ SystemTap ○○○○○○○○○ eBPF : bcc-tools, bpfftrace ○○○○○○○○○○○○○○○○○○○

```
crash> struct -o ptlrpc_request
```

```
struct ptlrpc_request {
    [0x0] int rq_type;
    [0x4] int rq_status;
    ...
    [0x48] struct lustre_msg_v2 *rq_reqmsg;
    [0x50] struct lustre_msg_v2 *rq_repmsg;
    [0x58] __u64 rq_transno;
    [0x60] __u64 rq_xid;
    [0x68] __u64 rq_mbits;
    ...
    [0x390] struct obd_export *rq_export;
    [0x398] struct obd_import *rq_import;
    [0x3a0] lnet_nid_t rq_self;
    [0x3a8] struct lnet_process_id rq_peer;
    [0x3b8] struct lnet_process_id rq_source;
    [0x3c8] time_t rq_timeout;
    [0x3d0] time64_t rq_sent;
    [0x3d8] time64_t rq_deadline;
    [0x3e0] struct req_capsule rq_pill;
}
SIZE: 0x450
```

- 0x450 = 1104 : kmalloc-2048 ou ptlrpc_cache
- à part le nid dans `rq_peer`, le reste est assez difficilement exploitable

Intro
○○○○○

Crash
○○○●○○○○○

Perf
○○○○○○○○○○○

SystemTap
○○○○○○○○○

eBPF : bcc-tools, bpfftrace
○○○○○○○○○○○○○○

```
crash> struct ptlrpc_request ffff917cb168a050
...
rq_export = 0xffff917c96f59000,
rq_import = 0x0,
rq_self = 0x200000ac80001,
rq_peer = {
    nid = 0x200000ac80002,
    pid = 0x3039
},
rq_source = {
    nid = 0x200000ac80002,
    pid = 0x3039
},
rq_timeout = 0x0,
rq_sent = 0x0,
...
```

Intro



Crash



Perf



SystemTap



eBPF : bcc-tools, bpfftrace



- 1s bloqués dans un dossier utilisateur pendant que je préparais mes slides...
- Petit tour rapide : vous devriez pouvoir refaire le même type de raisonnements

Intro

○○○○○

Crash

○○●○○○

Perf

○○○○○○○○○○

SystemTap

○○○○○○○○

eBPF : bcc-tools, bpftrace

○○○○○○○○○○○○○○

```

-----
"mdt00_069"/18640,"mdt00_100"/18689,
"mdt00_117"/18714,"mdt01_041"/25805,
"mdt01_058"/18667,"mdt01_088"/18734,
"mdt01_109"/7506 (7)
-----

```

```

#0 __schedule ffffffff97768a32
#1 schedule ffffffff97768ed9
#2 ldlm_completion_ast ffffffff0ea9aad
#3 ldlm_cli_enqueue_local ffffffff0eaab63
#4 mdt_object_local_lock ffffffff14d1dc2
#5 mdt_object_lock_internal ffffffff14d240e
#6 mdt_getattr_name_lock ffffffff14d31b4
#7 mdt_intent_getattr ffffffff14d44d0
#8 mdt_intent_opc ffffffff14d04bb
#9 mdt_intent_policy ffffffff14d8d58
#10 ldlm_lock_enqueue ffffffff0e8f2dd
#11 ldlm_handle_enqueue0 ffffffff0eb8c03
#12 tgt_enqueue ffffffff0f3e892
#13 tgt_request_handle ffffffff0f427ca
#14 ptlrpc_server_handle_request ffffffff0eeb05b
#15 ptlrpc_main ffffffff0eee7a2
#16 kthread ffffffff970c1da1

```

```

-----
"mdt01_033"/25786
-----

```

```

#0 __schedule ffffffff97768a32
#1 schedule ffffffff97768ed9
#2 ldlm_completion_ast ffffffff0ea9aad
#3 ldlm_cli_enqueue_local ffffffff0eaab63
#4 mdt_object_local_lock ffffffff14d1dc2
#5 mdt_object_lock_internal ffffffff14d240e
#6 mdt_reint_object_lock ffffffff14d26f0
#7 mdt_reint_unlink ffffffff14e84b2
#8 mdt_reint_rec ffffffff14ebc53
#9 mdt_reint_internal ffffffff14cd18b
#10 mdt_reint ffffffff14d8fa7
#11 tgt_request_handle ffffffff0f427ca
#12 ptlrpc_server_handle_request ffffffff0eeb05b
#13 ptlrpc_main ffffffff0eee7a2
#14 kthread ffffffff970c1da1

```

Intro Crash Perf SystemTap eBPF : bcc-tools, bpfftrace

○○○○○ ○○●○○○ ○○○○○○○○○○ ○○○○○○○○○○ ○○○○○○○○○○○○○○

```

crash> ps -l
[3621688788852009] [IN]  PID: 28644  TASK: ffff94a47bf84100  CPU: 2   COMMAND: "mdt00_013"
... 1700s ago
[3619975098195335] [IN]  PID: 6743   TASK: ffff94813e660000  CPU: 7   COMMAND: "kworker/7:1"
[3619964995355868] [IN]  PID: 18640  TASK: ffff9479eabb30c0  CPU: 0   COMMAND: "mdt00_069"
[3619922465196662] [IN]  PID: 24127  TASK: ffff94a47e22a080  CPU: 14  COMMAND: "mdt_seqm_0002"
... 6336s ago
[3615352056840100] [IN]  PID: 29072  TASK: ffff94a39afa6180  CPU: 15  COMMAND: "kworker/15:0"
[3615187577148168] [IN]  PID: 18714  TASK: ffff94a4785f9040  CPU: 0   COMMAND: "mdt00_117"
[3615157566040146] [IN]  PID: 18689  TASK: ffff94847fd6b0c0  CPU: 3   COMMAND: "mdt00_100"
[3614626563207010] [IN]  PID: 18667  TASK: ffff94a47e5f5140  CPU: 12  COMMAND: "mdt01_058"
[3614100469065774] [IN]  PID: 27971  TASK: ffff948472a8e180  CPU: 10  COMMAND: "cfs_rh_03"
[3614025495970123] [IN]  PID: 25805  TASK: ffff94847d738000  CPU: 8   COMMAND: "mdt01_041"
[3612352026539768] [IN]  PID: 3482   TASK: ffff94a47ebc0000  CPU: 8   COMMAND: "kworker/8:1"
[3612142134798789] [IN]  PID: 30482  TASK: ffff949e91689040  CPU: 11  COMMAND: "kworker/11:0"
[3612035313324469] [IN]  PID: 7506   TASK: ffff94a3a26cb0c0  CPU: 12  COMMAND: "mdt01_109"
[3611323248826132] [IN]  PID: 18734  TASK: ffff94a3e3e81040  CPU: 11  COMMAND: "mdt01_088"
[3611260896608780] [IN]  PID: 21311  TASK: ffff949ecd242080  CPU: 8   COMMAND: "kworker/u34:1"
[3608862966498563] [IN]  PID: 25786  TASK: ffff947dadf61040  CPU: 12  COMMAND: "mdt01_033"
[3605752338202911] [IN]  PID: 32477  TASK: ffff949f242c6180  CPU: 12  COMMAND: "kworker/12:1"
...

```

- stuck quelques heures !
- unlink en premier

- Il reste plein de commandes pas explorées ici
 - dis(assemble), kmem, mount, dev, net, rd, bpf, etc. . .
- Vraiment lire le détail des backtraces
 - En réfléchissant un peu on peut retrouver exactement quelle position contient quelle variable sans deviner avec les slabs, mais ça serait une autre formation. . .
- extensions crash :
 - extension lustre pour tirer un dk du dump ou autre
 - <https://people.redhat.com/anderson/extensions.html>

1 Intro

2 Crash

3 Perf

- perf probe
- flame graph

4 SystemTap

5 eBPF : bcc-tools, bpftrace

Intro



Crash



Perf



SystemTap



eBPF : bcc-tools, bpftrace



- Introduire des tracepoints un peu n'importe où dans le kernel (ou l'userspace !)

Intro



Crash



Perf



SystemTap



eBPF : bcc-tools, bpftrace



- Pas le temps de présenter en détail
- Allez voir le site de Brendan Gregg
 - cpu flame graphs¹
 - off-cpu flame graphs²

1. <http://www.brendangregg.com/FlameGraphs/cpuflamegraphs.html>

2. <http://www.brendangregg.com/FlameGraphs/offcpuflamegraphs.html>

- 1 Intro
- 2 Crash
- 3 Perf
- 4 SystemTap**
- 5 eBPF : bcc-tools, bpftrace

- Plus puissant / intrusif que perf
- Compile un petit module kernel et l'insert : on peut tout casser
 - Mais du coup on peut l'utiliser pour contourner des problèmes !

```
#!/usr/bin/stap -g

global EIO = -5
global EAGAIN = -11

probe begin {
    print("LU-5642 stap started. Press ^C to exit\n")
}

function syslog(msg: string) %{
    printk("stap lu_5642: %s\n", STAP_ARG_msg);
}%

probe module("lustre").function("ll_xattr_cache_refill").return {
    if ($return == EIO && uid() > 1000) {
        slurm_jobid = env_var("SLURM_JOBID")
        syslog(sprintf("JOBID %s: %s(%d) got EIO, changing to EAGAIN"
            " (inode %p)",
            slurm_jobid, execname(), pid(), $inode))
        $return = EAGAIN
    }
}
```


- Plus puissant que perf pour dumper des types complexes : \$var\$\$
- Essayons de faire comme la probe perf...

Intro



Crash



Perf



SystemTap



eBPF : bcc-tools, bpftrace



- <https://sourceware.org/systemtap/langref.pdf>
- <https://github.com/calio/systemtap-cheat-sheet>

- 1 Intro
- 2 Crash
- 3 Perf
- 4 SystemTap
- 5 eBPF : bcc-tools, bpftrace
 - bcc-tools
 - bpftrace
 - bcc/bpftrace internals



Merci !

Intro



Crash



Perf



SystemTap



eBPF : bcc-tools, bpftrace

