

DE LA RECHERCHE À L'INDUSTRIE



www.cea.fr

Kernel debugging

Tools to understand whatever it is that is happening in there

Dominique Martinet

CEA

January 9, 2020

- 1 Introduction
 - Foreword
 - Hands on setup

2 Crash

3 Perf

4 SystemTap

5 eBPF: bcc-tools, bpftrace

1 Introduction

2 Crash

- lustre client LBUG
- Lustre server load
- Glance at another example

3 Perf

4 SystemTap

5 eBPF: bcc-tools, bpftrace


```
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"
```

- ### Crash notes
- First column (brackets) contains the last scheduled timestamp of each task (in nanosecond from boot)
 - Processes are sorted on this column: helps find group of stuck threads
 - 1834107714279 – 1833456118992 = 652ms
 - dmesg will show traces from 120s in UN or RU state there, this case isn't realistic
 - dmesg | tail -> [1834.107266] ...


```
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,
```

Crash notes

- `lr_name` often a fid, but on OST here it is an oid
- `lr_refcount` refcount (duh) associated with the resource
- `lr_granted/waiting` lists of locks

```
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",
```

Crash notes

- Going through multiple structures is a bit painful, but can be done

Introduction

○○○○

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

```

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

- ## Notes
- double-linked list
 - Often used with an independant: `list -H`
 - i.e. `lustre/ldlm/ldlm_internal.h:38: extern struct list_head ldlm_srv_namespace_list;`
 - `/!\` same type = sometimes hard to tell appart

```
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
```

Crash notes

- Here (read code!) lr_granted/lr_waiting are similar list “heads”
- Links to struct ldlm_lock on field 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
...

```

Crash notes

- -S is faster than -s, but no symbol resolution
- multiple LCK_PW granted but most are PR, should be incompatible but l_req_extents are different: work on independant ranges.

Introduction

○○○

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
  },
```

Crash notes

- nid = lnd type, lnd number, ip/id
 - Here: 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

```


1 Introduction

2 Crash

3 Perf

- perf probe
- flame graph

4 SystemTap

5 eBPF: bcc-tools, bpftrace

1 Introduction

2 Crash

3 Perf

4 SystemTap

5 eBPF: bcc-tools, bpftrace

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



Thanks!

Introduction



Crash



Perf



SystemTap



eBPF: bcc-tools, bpftrace

