

---

Subject: [PATCH v8 0/9] per-cgroup tcp memory pressure controls

Posted by [Glauber Costa](#) on Mon, 05 Dec 2011 21:34:54 GMT

[View Forum Message](#) <> [Reply to Message](#)

---

Hi,

This is my new attempt to fix all the concerns that were raised during the last iteration.

I should highlight:

- 1) proc information is kept intact. (although I kept the wrapper functions) it will be submitted as a follow up patch so it can get the attention it deserves
- 2) sockets now hold a reference to memcg. sockets can be alive even after the task is gone, so we don't bother with between cgroups movements.  
To be able to release resources more easily in this cenario, the parent pointer in struct cg\_proto was replaced by a memcg object. We then iterate through its pointer (which is cleaner anyway)

The rest should be mostly the same except for small fixes and style changes.

Glauber Costa (9):

Basic kernel memory functionality for the Memory Controller

foundations of per-cgroup memory pressure controlling.

socket: initial cgroup code.

tcp memory pressure controls

per-netns ipv4 sysctl\_tcp\_mem

tcp buffer limitation: per-cgroup limit

Display current tcp memory allocation in kmem cgroup

Display current tcp failcnt in kmem cgroup

Display maximum tcp memory allocation in kmem cgroup

```
Documentation/cgroups/memory.txt | 46 ++++++-
include/linux/memcontrol.h       | 23 +++++
include/net/netns/ipv4.h         | 1 +
include/net/sock.h               | 239 +++++
include/net/tcp.h                 | 4 +-
include/net/tcp_memcontrol.h     | 19 +++
init/Kconfig                     | 11 ++
mm/memcontrol.c                  | 189 +++++
net/core/sock.c                  | 118 +++++
net/ipv4/Makefile                | 1 +
net/ipv4/af_inet.c               | 2 +
net/ipv4/proc.c                  | 6 +-
net/ipv4/sysctl_net_ipv4.c       | 65 ++++++-
net/ipv4/tcp.c                   | 11 +--
net/ipv4/tcp_input.c             | 12 +-
net/ipv4/tcp_ipv4.c              | 14 +-

```

```
net/ipv4/tcp_memcontrol.c      | 272 ++++++
net/ipv4/tcp_output.c          | 2 +-
net/ipv4/tcp_timer.c           | 2 +-
net/ipv6/af_inet6.c            | 2 +
net/ipv6/tcp_ipv6.c            | 8 +-
21 files changed, 968 insertions(+), 79 deletions(-)
create mode 100644 include/net/tcp_memcontrol.h
create mode 100644 net/ipv4/tcp_memcontrol.c
```

--

1.7.6.4

---

Subject: [PATCH v8 1/9] Basic kernel memory functionality for the Memory Controller

Posted by [Glauber Costa](#) on Mon, 05 Dec 2011 21:34:55 GMT

[View Forum Message](#) <> [Reply to Message](#)

---

This patch lays down the foundation for the kernel memory component of the Memory Controller.

As of today, I am only laying down the following files:

- \* memory.independent\_kmem\_limit
- \* memory.kmem.limit\_in\_bytes (currently ignored)
- \* memory.kmem.usage\_in\_bytes (always zero)

Signed-off-by: Glauber Costa <glommer@parallels.com>

Reviewed-by: Kirill A. Shutemov <kirill@shutemov.name>

CC: Paul Menage <paul@paulmenage.org>

CC: Greg Thelen <gthelen@google.com>

---

```
Documentation/cgroups/memory.txt | 40 ++++++
init/Kconfig                      | 11 ++++
mm/memcontrol.c                   | 103 ++++++
3 files changed, 147 insertions(+), 7 deletions(-)
```

diff --git a/Documentation/cgroups/memory.txt b/Documentation/cgroups/memory.txt

index cc0ebc5..f245324 100644

--- a/Documentation/cgroups/memory.txt

+++ b/Documentation/cgroups/memory.txt

@@ -44,8 +44,9 @@ Features:

- oom-killer disable knob and oom-notifier
- Root cgroup has no limit controls.
- Kernel memory and Hugepages are not under control yet. We just manage
- pages on LRU. To add more controls, we have to take care of performance.
- + Hugepages is not under control yet. We just manage pages on LRU. To add more

+ controls, we have to take care of performance. Kernel memory support is work  
+ in progress, and the current version provides basically functionality.

Brief summary of control files.

@@ -56,8 +57,11 @@ Brief summary of control files.

(See 5.5 for details)

memory.memsw.usage\_in\_bytes # show current res\_counter usage for memory+Swap

(See 5.5 for details)

+ memory.kmem.usage\_in\_bytes # show current res\_counter usage for kmem only.

+ (See 2.7 for details)

memory.limit\_in\_bytes # set/show limit of memory usage

memory.memsw.limit\_in\_bytes # set/show limit of memory+Swap usage

+ memory.kmem.limit\_in\_bytes # if allowed, set/show limit of kernel memory

memory.failcnt # show the number of memory usage hits limits

memory.memsw.failcnt # show the number of memory+Swap hits limits

memory.max\_usage\_in\_bytes # show max memory usage recorded

@@ -72,6 +76,9 @@ Brief summary of control files.

memory.oom\_control # set/show oom controls.

memory.numa\_stat # show the number of memory usage per numa node

+ memory.independent\_kmem\_limit # select whether or not kernel memory limits are

+ independent of user limits

+

## 1. History

The memory controller has a long history. A request for comments for the memory

@@ -255,6 +262,35 @@ When oom event notifier is registered, event will be delivered.

per-zone-per-cgroup LRU (cgroup's private LRU) is just guarded by

zone->lru\_lock, it has no lock of its own.

## +2.7 Kernel Memory Extension (CONFIG\_CGROUP\_MEM\_RES\_CTLR\_KMEM)

+

+With the Kernel memory extension, the Memory Controller is able to limit

+the amount of kernel memory used by the system. Kernel memory is fundamentally

+different than user memory, since it can't be swapped out, which makes it

+possible to DoS the system by consuming too much of this precious resource.

+

+Some kernel memory resources may be accounted and limited separately from the

+main "kmem" resource. For instance, a slab cache that is considered important

+enough to be limited separately may have its own knobs.

+

+Kernel memory limits are not imposed for the root cgroup. Usage for the root

+cgroup may or may not be accounted.

+

+Memory limits as specified by the standard Memory Controller may or may not

+take kernel memory into consideration. This is achieved through the file

+memory.independent\_kmem\_limit. A Value different than 0 will allow for kernel

- +memory to be controlled separately.
- +
- +When kernel memory limits are not independent, the limit values set in
- +memory.kmem files are ignored.
- +
- +Currently no soft limit is implemented for kernel memory. It is future work
- +to trigger slab reclaim when those limits are reached.

## +2.7.1 Current Kernel Memory resources accounted

+None

## 3. User Interface

### 0. Configuration

```
diff --git a/init/Kconfig b/init/Kconfig
```

```
index 43298f9..b8930d5 100644
```

```
--- a/init/Kconfig
```

```
+++ b/init/Kconfig
```

```
@ @ -689,6 +689,17 @ @ config CGROUP_MEM_RES_CTLR_SWAP_ENABLED
```

For those who want to have the feature enabled by default should  
select this option (if, for some reason, they need to disable it  
then swapaccount=0 does the trick).

```
+config CGROUP_MEM_RES_CTLR_KMEM
```

```
+ bool "Memory Resource Controller Kernel Memory accounting (EXPERIMENTAL)"
```

```
+ depends on CGROUP_MEM_RES_CTLR && EXPERIMENTAL
```

```
+ default n
```

```
+ help
```

+ The Kernel Memory extension for Memory Resource Controller can limit  
+ the amount of memory used by kernel objects in the system. Those are  
+ fundamentally different from the entities handled by the standard  
+ Memory Controller, which are page-based, and can be swapped. Users of  
+ the kmem extension can use it to guarantee that no group of processes  
+ will ever exhaust kernel resources alone.

```
config CGROUP_PERF
```

```
bool "Enable perf_event per-cpu per-container group (cgroup) monitoring"
```

```
diff --git a/mm/memcontrol.c b/mm/memcontrol.c
```

```
index 6aff93c..3becb24 100644
```

```
--- a/mm/memcontrol.c
```

```
+++ b/mm/memcontrol.c
```

```
@ @ -227,6 +227,10 @ @ struct mem_cgroup {
```

```
*/
```

```
struct res_counter memsw;
```

```
/*
```

```
+ * the counter to account for kmem usage.
```

```
+ */
```

```
+ struct res_counter kmem;
```

```

+ /*
+  * Per cgroup active and inactive list, similar to the
+  * per zone LRU lists.
+  */
@@ -277,6 +281,11 @@ struct mem_cgroup {
+ /*
+  unsigned long  move_charge_at_immigrate;
+  */
+  * Should kernel memory limits be stabilshed independently
+  * from user memory ?
+  */
+ int  kmem_independent_accounting;
+ /*
+  * percpu counter.
+  */
+  struct mem_cgroup_stat_cpu *stat;
@@ -344,9 +353,14 @@ enum charge_type {
};

/* for encoding cft->private value on file */
#define _MEM  (0)
#define _MEMSWAP  (1)
#define _OOM_TYPE  (2)
+
+enum mem_type {
+  _MEM = 0,
+  _MEMSWAP,
+  _OOM_TYPE,
+  _KMEM,
+};
+
+
#define MEMFILE_PRIVATE(x, val) (((x) << 16) | (val))
#define MEMFILE_TYPE(val) (((val) >> 16) & 0xffff)
#define MEMFILE_ATTR(val) ((val) & 0xffff)
@@ -3848,10 +3862,17 @@ static inline u64 mem_cgroup_usage(struct mem_cgroup *memcg,
bool swap)
    u64 val;

    if (!mem_cgroup_is_root(memcg)) {
+  val = 0;
+  #ifdef CONFIG_CGROUP_MEM_RES_CTLR_KMEM
+  if (!memcg->kmem_independent_accounting)
+  val = res_counter_read_u64(&memcg->kmem, RES_USAGE);
+  #endif
+  if (!swap)
-  return res_counter_read_u64(&memcg->res, RES_USAGE);
+  val += res_counter_read_u64(&memcg->res, RES_USAGE);
+  else

```

```

- return res_counter_read_u64(&memcg->memsw, RES_USAGE);
+ val += res_counter_read_u64(&memcg->memsw, RES_USAGE);
+
+ return val;
}

val = mem_cgroup_recursive_stat(memcg, MEM_CGROUP_STAT_CACHE);
@@ -3884,6 +3905,11 @@ static u64 mem_cgroup_read(struct cgroup *cont, struct cftype *cft)
    else
        val = res_counter_read_u64(&memcg->memsw, name);
    break;
+ #ifdef CONFIG_CGROUP_MEM_RES_CTLR_KMEM
+ case _KMEM:
+ val = res_counter_read_u64(&memcg->kmem, name);
+ break;
+ #endif
    default:
        BUG();
        break;
@@ -4612,6 +4638,67 @@ static int mem_control_numa_stat_open(struct inode *unused, struct
file *file)
}
#endif /* CONFIG_NUMA */

+ #ifdef CONFIG_CGROUP_MEM_RES_CTLR_KMEM
+ static u64 kmem_limit_independent_read(struct cgroup *cgroup, struct cftype *cft)
+ {
+ return mem_cgroup_from_cont(cgroup)->kmem_independent_accounting;
+ }
+
+ static int kmem_limit_independent_write(struct cgroup *cgroup, struct cftype *cft,
+ u64 val)
+ {
+ struct mem_cgroup *memcg = mem_cgroup_from_cont(cgroup);
+ struct mem_cgroup *parent = parent_mem_cgroup(memcg);
+
+ val = !!val;
+
+ if (parent && parent->use_hierarchy &&
+ (val != parent->kmem_independent_accounting))
+ return -EINVAL;
+ /*
+  * TODO: We need to handle the case in which we are doing
+  * independent kmem accounting as authorized by our parent,
+  * but then our parent changes its parameter.
+  */
+ cgroup_lock();
+ memcg->kmem_independent_accounting = val;

```



```

@@ -4935,6 +5023,7 @@ mem_cgroup_create(struct cgroup_subsys *ss, struct cgroup *cont)
} else {
    res_counter_init(&memcg->res, NULL);
    res_counter_init(&memcg->memsw, NULL);
+ res_counter_init(&memcg->kmem, NULL);
}
    memcg->last_scanned_child = 0;
    memcg->last_scanned_node = MAX_NUMNODES;
@@ -4978,6 +5067,10 @@ static int mem_cgroup_populate(struct cgroup_subsys *ss,

    if (!ret)
        ret = register_memsw_files(cont, ss);
+
+ if (!ret)
+ ret = register_kmem_files(cont, ss);
+
    return ret;
}

--
1.7.6.4

```

---

Subject: [PATCH v8 2/9] foundations of per-cgroup memory pressure controlling.  
 Posted by [Glauber Costa](#) on Mon, 05 Dec 2011 21:34:56 GMT  
[View Forum Message](#) <> [Reply to Message](#)

---

This patch replaces all uses of struct sock fields' memory\_pressure, memory\_allocated, sockets\_allocated, and sysctl\_mem to accessor macros. Those macros can either receive a socket argument, or a mem\_cgroup argument, depending on the context they live in.

Since we're only doing a macro wrapping here, no performance impact at all is expected in the case where we don't have cgroups disabled.

Signed-off-by: Glauber Costa <glommer@parallels.com>  
 CC: David S. Miller <davem@davemloft.net>  
 CC: Hiroyouki Kamezawa <kamezawa.hiroyu@jp.fujitsu.com>  
 CC: Eric W. Biederman <ebiederm@xmission.com>  
 CC: Eric Dumazet <eric.dumazet@gmail.com>

```

---
include/net/sock.h | 96 ++++++
include/net/tcp.h | 3 +-
net/core/sock.c | 59 ++++++
net/ipv4/proc.c | 6 +-
net/ipv4/tcp_input.c | 12 +---
net/ipv4/tcp_ipv4.c | 4 +-
net/ipv4/tcp_output.c | 2 +-

```



```
net/ipv4/tcp_timer.c | 2 +-
net/ipv6/tcp_ipv6.c | 2 +-
9 files changed, 145 insertions(+), 41 deletions(-)
```

```
diff --git a/include/net/sock.h b/include/net/sock.h
index abb6e0f..5f43fd9 100644
--- a/include/net/sock.h
+++ b/include/net/sock.h
@@ -53,6 +53,7 @@
#include <linux/security.h>
#include <linux/slab.h>
#include <linux/uaccess.h>
+#include <linux/memcontrol.h>

#include <linux/filter.h>
#include <linux/rculist_nulls.h>
@@ -863,6 +864,99 @@ static inline void sk_refcnt_debug_release(const struct sock *sk)
#define sk_refcnt_debug_release(sk) do { } while (0)
#endif /* SOCK_REFCNT_DEBUG */

+static inline bool sk_has_memory_pressure(const struct sock *sk)
+{
+ return sk->sk_prot->memory_pressure != NULL;
+}
+
+static inline bool sk_under_memory_pressure(const struct sock *sk)
+{
+ if (!sk->sk_prot->memory_pressure)
+ return false;
+ return !!*sk->sk_prot->memory_pressure;
+}
+
+static inline void sk_leave_memory_pressure(struct sock *sk)
+{
+ int *memory_pressure = sk->sk_prot->memory_pressure;
+
+ if (memory_pressure && *memory_pressure)
+ *memory_pressure = 0;
+}
+
+static inline void sk_enter_memory_pressure(struct sock *sk)
+{
+ if (sk->sk_prot->enter_memory_pressure)
+ sk->sk_prot->enter_memory_pressure(sk);
+}
+
+static inline long sk_prot_mem_limits(const struct sock *sk, int index)
+{
```

```

+ long *prot = sk->sk_prot->sysctl_mem;
+ return prot[index];
+}
+
+static inline long
+sk_memory_allocated(const struct sock *sk)
+{
+ struct proto *prot = sk->sk_prot;
+ return atomic_long_read(prot->memory_allocated);
+}
+
+static inline long
+sk_memory_allocated_add(struct sock *sk, int amt)
+{
+ struct proto *prot = sk->sk_prot;
+ return atomic_long_add_return(amt, prot->memory_allocated);
+}
+
+static inline void
+sk_memory_allocated_sub(struct sock *sk, int amt)
+{
+ struct proto *prot = sk->sk_prot;
+ atomic_long_sub(amt, prot->memory_allocated);
+}
+
+static inline void sk_sockets_allocated_dec(struct sock *sk)
+{
+ struct proto *prot = sk->sk_prot;
+ percpu_counter_dec(prot->sockets_allocated);
+}
+
+static inline void sk_sockets_allocated_inc(struct sock *sk)
+{
+ struct proto *prot = sk->sk_prot;
+ percpu_counter_inc(prot->sockets_allocated);
+}
+
+static inline int
+sk_sockets_allocated_read_positive(struct sock *sk)
+{
+ struct proto *prot = sk->sk_prot;
+
+ return percpu_counter_sum_positive(prot->sockets_allocated);
+}
+
+static inline int
+proto_sockets_allocated_sum_positive(struct proto *prot)
+{

```

```

+ return percpu_counter_sum_positive(prot->sockets_allocated);
+}
+
+static inline long
+proto_memory_allocated(struct proto *prot)
+{
+ return atomic_long_read(prot->memory_allocated);
+}
+
+static inline bool
+proto_memory_pressure(struct proto *prot)
+{
+ if (!prot->memory_pressure)
+ return false;
+ return !!*prot->memory_pressure;
+}
+

#ifdef CONFIG_PROC_FS
/* Called with local bh disabled */
@@ -1670,7 +1764,7 @@ static inline struct page *sk_stream_alloc_page(struct sock *sk)

    page = alloc_pages(sk->sk_allocation, 0);
    if (!page) {
- sk->sk_prot->enter_memory_pressure(sk);
+ sk_enter_memory_pressure(sk);
    sk_stream_moderate_sndbuf(sk);
    }
    return page;
diff --git a/include/net/tcp.h b/include/net/tcp.h
index bb18c4d..f080e0b 100644
--- a/include/net/tcp.h
+++ b/include/net/tcp.h
@@ -44,6 +44,7 @@
#include <net/dst.h>

#include <linux/seq_file.h>
+#include <linux/memcontrol.h>

extern struct inet_hashinfo tcp_hashinfo;

@@ -285,7 +286,7 @@ static inline bool tcp_too_many_orphans(struct sock *sk, int shift)
}

if (sk->sk_wmem_queued > SOCK_MIN_SNDBUF &&
- atomic_long_read(&tcp_memory_allocated) > sysctl_tcp_mem[2])
+ sk_memory_allocated(sk) > sk_prot_mem_limits(sk, 2))
    return true;

```

```

    return false;
}
diff --git a/net/core/sock.c b/net/core/sock.c
index 4ed7b1d..2b86d24 100644
--- a/net/core/sock.c
+++ b/net/core/sock.c
@@ -1288,7 +1288,7 @@ struct sock *sk_clone(const struct sock *sk, const gfp_t priority)
    newsk->sk_wq = NULL;

    if (newsk->sk_prot->sockets_allocated)
-   percpu_counter_inc(newsk->sk_prot->sockets_allocated);
+   sk_sockets_allocated_inc(newsk);

    if (sock_flag(newsk, SOCK_TIMESTAMP) ||
        sock_flag(newsk, SOCK_TIMESTAMPING_RX_SOFTWARE))
@@ -1679,28 +1679,28 @@ int __sk_mem_schedule(struct sock *sk, int size, int kind)
    long allocated;

    sk->sk_forward_alloc += amt * SK_MEM_QUANTUM;
-   allocated = atomic_long_add_return(amt, prot->memory_allocated);
+
+   allocated = sk_memory_allocated_add(sk, amt);

    /* Under limit. */
-   if (allocated <= prot->sysctl_mem[0]) {
-       if (prot->memory_pressure && *prot->memory_pressure)
-           *prot->memory_pressure = 0;
-       return 1;
-   }
+   if (allocated <= sk_prot_mem_limits(sk, 0))
+       sk_leave_memory_pressure(sk);

    /* Under pressure. */
-   if (allocated > prot->sysctl_mem[1])
-       if (prot->enter_memory_pressure)
-           prot->enter_memory_pressure(sk);
+   if (allocated > sk_prot_mem_limits(sk, 1))
+       sk_enter_memory_pressure(sk);

    /* Over hard limit. */
-   if (allocated > prot->sysctl_mem[2])
+   if (allocated > sk_prot_mem_limits(sk, 2))
        goto suppress_allocation;

    /* guarantee minimum buffer size under pressure */
    if (kind == SK_MEM_RECV) {
        if (atomic_read(&sk->sk_rmem_alloc) < prot->sysctl_rmem[0])
            return 1;
    }

```

```

+
} else { /* SK_MEM_SEND */
    if (sk->sk_type == SOCK_STREAM) {
        if (sk->sk_wmem_queued < prot->sysctl_wmem[0])
@@ -1710,13 +1708,13 @@ int __sk_mem_schedule(struct sock *sk, int size, int kind)
        return 1;
    }

- if (prot->memory_pressure) {
+ if (sk_has_memory_pressure(sk)) {
    int alloc;

- if (!*prot->memory_pressure)
+ if (!sk_under_memory_pressure(sk))
        return 1;
- alloc = percpu_counter_read_positive(prot->sockets_allocated);
- if (prot->sysctl_mem[2] > alloc *
+ alloc = sk_sockets_allocated_read_positive(sk);
+ if (sk_prot_mem_limits(sk, 2) > alloc *
        sk_mem_pages(sk->sk_wmem_queued +
        atomic_read(&sk->sk_rmem_alloc) +
        sk->sk_forward_alloc))
@@ -1739,7 +1737,9 @@ suppress_allocation:

    /* Alas. Undo changes. */
    sk->sk_forward_alloc -= amt * SK_MEM_QUANTUM;
- atomic_long_sub(amt, prot->memory_allocated);
+
+ sk_memory_allocated_sub(sk, amt);
+
    return 0;
}
EXPORT_SYMBOL(__sk_mem_schedule);
@@ -1750,15 +1750,13 @@ EXPORT_SYMBOL(__sk_mem_schedule);
*/
void __sk_mem_reclaim(struct sock *sk)
{
- struct proto *prot = sk->sk_prot;
-
- atomic_long_sub(sk->sk_forward_alloc >> SK_MEM_QUANTUM_SHIFT,
- prot->memory_allocated);
+ sk_memory_allocated_sub(sk,
+ sk->sk_forward_alloc >> SK_MEM_QUANTUM_SHIFT);
    sk->sk_forward_alloc &= SK_MEM_QUANTUM - 1;

- if (prot->memory_pressure && *prot->memory_pressure &&
-     (atomic_long_read(prot->memory_allocated) < prot->sysctl_mem[0]))
-     *prot->memory_pressure = 0;

```

```

+ if (sk_under_memory_pressure(sk) &&
+     (sk_memory_allocated(sk) < sk_prot_mem_limits(sk, 0)))
+ sk_leave_memory_pressure(sk);
+ }
EXPORT_SYMBOL(__sk_mem_reclaim);

@@ -2474,16 +2472,27 @@ static char proto_method_implemented(const void *method)
{
    return method == NULL ? 'n' : 'y';
}
+static long sock_prot_memory_allocated(struct proto *proto)
+{
+ return proto->memory_allocated != NULL ? proto_memory_allocated(proto): -1L;
+}
+
+static char *sock_prot_memory_pressure(struct proto *proto)
+{
+ return proto->memory_pressure != NULL ?
+ proto_memory_pressure(proto) ? "yes" : "no" : "NI";
+}

static void proto_seq_printf(struct seq_file *seq, struct proto *proto)
{
+
    seq_printf(seq, "%-9s %4u %6d %6ld  %-3s %6u  %-3s %-10s "
        "%2c %2c %2c %2c %2c %2c %2c %2c %2c %2c %2c %2c %2c %2c %2c %2c %2c %2c\n",
        proto->name,
        proto->obj_size,
        sock_prot_inuse_get(seq_file_net(seq), proto),
-   proto->memory_allocated != NULL ? atomic_long_read(proto->memory_allocated) : -1L,
-   proto->memory_pressure != NULL ? *proto->memory_pressure ? "yes" : "no" : "NI",
+   sock_prot_memory_allocated(proto),
+   sock_prot_memory_pressure(proto),
        proto->max_header,
        proto->slab == NULL ? "no" : "yes",
        module_name(proto->owner),
diff --git a/net/ipv4/proc.c b/net/ipv4/proc.c
index 466ea8b..91be152 100644
--- a/net/ipv4/proc.c
+++ b/net/ipv4/proc.c
@@ -56,17 +56,17 @@ static int sockstat_seq_show(struct seq_file *seq, void *v)

    local_bh_disable();
    orphans = percpu_counter_sum_positive(&tcp_orphan_count);
- sockets = percpu_counter_sum_positive(&tcp_sockets_allocated);
+ sockets = proto_sockets_allocated_sum_positive(&tcp_prot);
    local_bh_enable();

```

```

socket_seq_show(seq);
seq_printf(seq, "TCP: inuse %d orphan %d tw %d alloc %d mem %ld\n",
    sock_prot_inuse_get(net, &tcp_prot), orphans,
    tcp_death_row.tw_count, sockets,
-   atomic_long_read(&tcp_memory_allocated));
+   proto_memory_allocated(&tcp_prot));
seq_printf(seq, "UDP: inuse %d mem %ld\n",
    sock_prot_inuse_get(net, &udp_prot),
-   atomic_long_read(&udp_memory_allocated));
+   proto_memory_allocated(&udp_prot));
seq_printf(seq, "UDPLITE: inuse %d\n",
    sock_prot_inuse_get(net, &udplite_prot));
seq_printf(seq, "RAW: inuse %d\n",
diff --git a/net/ipv4/tcp_input.c b/net/ipv4/tcp_input.c
index 52b5c2d..b64b5e8 100644
--- a/net/ipv4/tcp_input.c
+++ b/net/ipv4/tcp_input.c
@@ -322,7 +322,7 @@ static void tcp_grow_window(struct sock *sk, const struct sk_buff *skb)
/* Check #1 */
if (tp->rcv_ssthresh < tp->window_clamp &&
    (int)tp->rcv_ssthresh < tcp_space(sk) &&
-   !tcp_memory_pressure) {
+   !sk_under_memory_pressure(sk)) {
    int incr;

/* Check #2. Increase window, if skb with such overhead
@@ -411,8 +411,8 @@ static void tcp_clamp_window(struct sock *sk)

if (sk->sk_rcvbuf < sysctl_tcp_rmem[2] &&
    !(sk->sk_userlocks & SOCK_RCVBUF_LOCK) &&
-   !tcp_memory_pressure &&
-   atomic_long_read(&tcp_memory_allocated) < sysctl_tcp_mem[0]) {
+   !sk_under_memory_pressure(sk) &&
+   sk_memory_allocated(sk) < sk_prot_mem_limits(sk, 0)) {
    sk->sk_rcvbuf = min(atomic_read(&sk->sk_rmem_alloc),
        sysctl_tcp_rmem[2]);
}
@@ -486,7 +486,7 @@ static int tcp_prune_queue(struct sock *sk)

if (atomic_read(&sk->sk_rmem_alloc) >= sk->sk_rcvbuf)
    tcp_clamp_window(sk);
- else if (tcp_memory_pressure)
+ else if (sk_under_memory_pressure(sk))
    tp->rcv_ssthresh = min(tp->rcv_ssthresh, 4U * tp->advmss);

tcp_collapse_ofo_queue(sk);
@@ -493,11 +493,11 @@ static int tcp_should_expand_sndbuf(const struct sock *sk)

```

```

return 0;

/* If we are under global TCP memory pressure, do not expand. */
- if (tcp_memory_pressure)
+ if (sk_under_memory_pressure(sk))
    return 0;

/* If we are under soft global TCP memory pressure, do not expand. */
- if (atomic_long_read(&tcp_memory_allocated) >= sysctl_tcp_mem[0])
+ if (sk_memory_allocated(sk) >= sk_prot_mem_limits(sk, 0))
    return 0;

/* If we filled the congestion window, do not expand. */
diff --git a/net/ipv4/tcp_ipv4.c b/net/ipv4/tcp_ipv4.c
index a744315..d1f4bf8 100644
--- a/net/ipv4/tcp_ipv4.c
+++ b/net/ipv4/tcp_ipv4.c
@@ -1915,7 +1915,7 @@ static int tcp_v4_init_sock(struct sock *sk)
    sk->sk_rcvbuf = sysctl_tcp_rmem[1];

    local_bh_disable();
- percpu_counter_inc(&tcp_sockets_allocated);
+ sk_sockets_allocated_inc(sk);
    local_bh_enable();

    return 0;
@@ -1971,7 +1971,7 @@ void tcp_v4_destroy_sock(struct sock *sk)
    tp->cookie_values = NULL;
}

- percpu_counter_dec(&tcp_sockets_allocated);
+ sk_sockets_allocated_dec(sk);
}
EXPORT_SYMBOL(tcp_v4_destroy_sock);

diff --git a/net/ipv4/tcp_output.c b/net/ipv4/tcp_output.c
index 980b98f..b378490 100644
--- a/net/ipv4/tcp_output.c
+++ b/net/ipv4/tcp_output.c
@@ -1919,7 +1919,7 @@ u32 __tcp_select_window(struct sock *sk)
    if (free_space < (full_space >> 1)) {
        icsk->icsk_ack.quick = 0;

- if (tcp_memory_pressure)
+ if (sk_under_memory_pressure(sk))
        tp->rcv_ssthresh = min(tp->rcv_ssthresh,
                                4U * tp->advmss);

```



```

diff --git a/net/ipv4/tcp_timer.c b/net/ipv4/tcp_timer.c
index 2e0f0af..d6ddacb 100644
--- a/net/ipv4/tcp_timer.c
+++ b/net/ipv4/tcp_timer.c
@@ -261,7 +261,7 @@ static void tcp_delack_timer(unsigned long data)
 }

out:
- if (tcp_memory_pressure)
+ if (sk_under_memory_pressure(sk))
    sk_mem_reclaim(sk);
out_unlock:
    bh_unlock_sock(sk);
diff --git a/net/ipv6/tcp_ipv6.c b/net/ipv6/tcp_ipv6.c
index 36131d1..e666768 100644
--- a/net/ipv6/tcp_ipv6.c
+++ b/net/ipv6/tcp_ipv6.c
@@ -1995,7 +1995,7 @@ static int tcp_v6_init_sock(struct sock *sk)
    sk->sk_rcvbuf = sysctl_tcp_rmem[1];

    local_bh_disable();
- percpu_counter_inc(&tcp_sockets_allocated);
+ sk_sockets_allocated_inc(sk);
    local_bh_enable();

    return 0;
--
1.7.6.4

```

---

Subject: [PATCH v8 3/9] socket: initial cgroup code.  
 Posted by [Glauber Costa](#) on Mon, 05 Dec 2011 21:34:57 GMT  
[View Forum Message](#) <> [Reply to Message](#)

---

The goal of this work is to move the memory pressure tcp controls to a cgroup, instead of just relying on global conditions.

To avoid excessive overhead in the network fast paths, the code that accounts allocated memory to a cgroup is hidden inside a static\_branch(). This branch is patched out until the first non-root cgroup is created. So when nobody is using cgroups, even if it is mounted, no significant performance penalty should be seen.

This patch handles the generic part of the code, and has nothing tcp-specific.

Signed-off-by: Glauber Costa <glommer@parallels.com>  
CC: Kirill A. Shutemov <kirill@shutemov.name>  
CC: KAMEZAWA Hiroyuki <kamezawa.hiroyu@jp.fujitsu.com>  
CC: David S. Miller <davem@davemloft.net>  
CC: Eric W. Biederman <ebiederm@xmission.com>  
CC: Eric Dumazet <eric.dumazet@gmail.com>

---

```
Documentation/cgroups/memory.txt | 4 +-
include/linux/memcontrol.h       | 22 ++++++
include/net/sock.h                | 151 ++++++
mm/memcontrol.c                  | 46 ++++++
net/core/sock.c                  | 24 +++++-
5 files changed, 230 insertions(+), 17 deletions(-)
```

diff --git a/Documentation/cgroups/memory.txt b/Documentation/cgroups/memory.txt  
index f245324..23a8dc5 100644

```
--- a/Documentation/cgroups/memory.txt
+++ b/Documentation/cgroups/memory.txt
@@ -289,7 +289,9 @@ to trigger slab reclaim when those limits are reached.
```

### 2.7.1 Current Kernel Memory resources accounted

-None  
+\* sockets memory pressure: some sockets protocols have memory pressure  
+thresholds. The Memory Controller allows them to be controlled individually  
+per cgroup, instead of globally.

### 3. User Interface

```
diff --git a/include/linux/memcontrol.h b/include/linux/memcontrol.h
index b87068a..f15021b 100644
--- a/include/linux/memcontrol.h
+++ b/include/linux/memcontrol.h
@@ -85,6 +85,8 @@ extern struct mem_cgroup *try_get_mem_cgroup_from_page(struct page
*page);
extern struct mem_cgroup *mem_cgroup_from_task(struct task_struct *p);
extern struct mem_cgroup *try_get_mem_cgroup_from_mm(struct mm_struct *mm);

+extern struct mem_cgroup *parent_mem_cgroup(struct mem_cgroup *memcg);
+
static inline
int mm_match_cgroup(const struct mm_struct *mm, const struct mem_cgroup *cgroup)
{
@@ -381,5 +383,25 @@ mem_cgroup_print_bad_page(struct page *page)
}
#endif

+#ifdef CONFIG_INET
```

```

+enum {
+ UNDER_LIMIT,
+ SOFT_LIMIT,
+ OVER_LIMIT,
+};
+
+struct sock;
+#ifdef CONFIG_CGROUP_MEM_RES_CTLR_KMEM
+void sock_update_memcg(struct sock *sk);
+void sock_release_memcg(struct sock *sk);
+#else
+static inline void sock_update_memcg(struct sock *sk)
+{
+}
+static inline void sock_release_memcg(struct sock *sk)
+{
+}
+#endif /* CONFIG_CGROUP_MEM_RES_CTLR_KMEM */
+#endif /* CONFIG_INET */
#endif /* _LINUX_MEMCONTROL_H */

```

```
diff --git a/include/net/sock.h b/include/net/sock.h
```

```
index 5f43fd9..b8a63f8 100644
```

```
--- a/include/net/sock.h
```

```
+++ b/include/net/sock.h
```

```
@@ -54,6 +54,7 @@
```

```
#include <linux/slab.h>
```

```
#include <linux/uaccess.h>
```

```
#include <linux/memcontrol.h>
```

```
+#include <linux/res_counter.h>
```

```
#include <linux/filter.h>
```

```
#include <linux/rculist_nulls.h>
```

```
@@ -168,6 +169,7 @@ struct sock_common {
```

```
/* public: */
```

```
};
```

```
+struct cg_proto;
```

```
/**
```

```
 * struct sock - network layer representation of sockets
```

```
 * @__sk_common: shared layout with inet_timewait_sock
```

```
@@ -228,6 +230,7 @@ struct sock_common {
```

```
 * @sk_security: used by security modules
```

```
 * @sk_mark: generic packet mark
```

```
 * @sk_classid: this socket's cgroup classid
```

```
+ * @sk_cgrp: this socket's cgroup-specific proto data
```

```
 * @sk_write_pending: a write to stream socket waits to start
```

```
 * @sk_state_change: callback to indicate change in the state of the sock
```

```

* @sk_data_ready: callback to indicate there is data to be processed
@@ -339,6 +342,7 @@ struct sock {
#endif
__u32 sk_mark;
u32 sk_classid;
+ struct cg_proto *sk_cgrp;
void (*sk_state_change)(struct sock *sk);
void (*sk_data_ready)(struct sock *sk, int bytes);
void (*sk_write_space)(struct sock *sk);
@@ -834,6 +838,37 @@ struct proto {
#ifdef SOCK_REFCNT_DEBUG
atomic_t socks;
#endif
+#ifdef CONFIG_CGROUP_MEM_RES_CTLR_KMEM
+ /*
+  * cgroup specific init/deinit functions. Called once for all
+  * protocols that implement it, from cgroups populate function.
+  * This function has to setup any files the protocol want to
+  * appear in the kmem cgroup filesystem.
+  */
+ int (*init_cgroup)(struct cgroup *cgrp,
+ struct cgroup_subsys *ss);
+ void (*destroy_cgroup)(struct cgroup *cgrp,
+ struct cgroup_subsys *ss);
+ struct cg_proto *(*proto_cgroup)(struct mem_cgroup *memcg);
+#endif
+};
+
+struct cg_proto {
+ void (*enter_memory_pressure)(struct sock *sk);
+ struct res_counter *memory_allocated; /* Current allocated memory. */
+ struct percpu_counter *sockets_allocated; /* Current number of sockets. */
+ int *memory_pressure;
+ long *sysctl_mem;
+ /*
+  * memcg field is used to find which memcg we belong directly
+  * Each memcg struct can hold more than one cg_proto, so container_of
+  * won't really cut.
+  *
+  * The elegant solution would be having an inverse function to
+  * proto_cgroup in struct proto, but that means polluting the structure
+  * for everybody, instead of just for memcg users.
+  */
+ struct mem_cgroup *memcg;
+};

extern int proto_register(struct proto *prot, int alloc_slab);
@@ -852,7 +887,7 @@ static inline void sk_refcnt_debug_dec(struct sock *sk)

```

```

    sk->sk_prot->name, sk, atomic_read(&sk->sk_prot->socks));
}

-static inline void sk_refcnt_debug_release(const struct sock *sk)
+inline void sk_refcnt_debug_release(const struct sock *sk)
{
    if (atomic_read(&sk->sk_refcnt) != 1)
        printk(KERN_DEBUG "Destruction of the %s socket %p delayed, refcnt=%d\n",
@@ -864,6 +899,19 @@ static inline void sk_refcnt_debug_release(const struct sock *sk)
#define sk_refcnt_debug_release(sk) do { } while (0)
#endif /* SOCK_REFCNT_DEBUG */

#ifdef CONFIG_CGROUP_MEM_RES_CTLR_KMEM
+extern struct jump_label_key memcg_socket_limit_enabled;
#define mem_cgroup_sockets_enabled static_branch(&memcg_socket_limit_enabled)
#else
#define mem_cgroup_sockets_enabled 0
#endif
+
+static inline struct cg_proto *parent_cg_proto(struct proto *proto,
+        struct cg_proto *cg_proto)
+{
+    return proto->proto_cgroup(parent_mem_cgroup(cg_proto->memcg));
+}
+
static inline bool sk_has_memory_pressure(const struct sock *sk)
{
    return sk->sk_prot->memory_pressure != NULL;
@@ -873,6 +921,10 @@ static inline bool sk_under_memory_pressure(const struct sock *sk)
{
    if (!sk->sk_prot->memory_pressure)
        return false;
+
+    if (mem_cgroup_sockets_enabled && sk->sk_cgrp)
+        return !!*sk->sk_cgrp->memory_pressure;
+
    return !!*sk->sk_prot->memory_pressure;
}

@@ -880,52 +932,136 @@ static inline void sk_leave_memory_pressure(struct sock *sk)
{
    int *memory_pressure = sk->sk_prot->memory_pressure;

-    if (memory_pressure && *memory_pressure)
+    if (!memory_pressure)
+        return;
+
+    if (*memory_pressure)

```

```

    *memory_pressure = 0;
+
+ if (mem_cgroup_sockets_enabled && sk->sk_cgrp) {
+ struct cg_proto *cg_proto = sk->sk_cgrp;
+ struct proto *prot = sk->sk_prot;
+
+ for (; cg_proto; cg_proto = parent_cg_proto(prot, cg_proto))
+ if (*cg_proto->memory_pressure)
+ *cg_proto->memory_pressure = 0;
+ }
+
+ }

static inline void sk_enter_memory_pressure(struct sock *sk)
{
- if (sk->sk_prot->enter_memory_pressure)
- sk->sk_prot->enter_memory_pressure(sk);
+ if (!sk->sk_prot->enter_memory_pressure)
+ return;
+
+ if (mem_cgroup_sockets_enabled && sk->sk_cgrp) {
+ struct cg_proto *cg_proto = sk->sk_cgrp;
+ struct proto *prot = sk->sk_prot;
+
+ for (; cg_proto; cg_proto = parent_cg_proto(prot, cg_proto))
+ cg_proto->enter_memory_pressure(sk);
+ }
+
+ sk->sk_prot->enter_memory_pressure(sk);
+ }

static inline long sk_prot_mem_limits(const struct sock *sk, int index)
{
    long *prot = sk->sk_prot->sysctl_mem;
+ if (mem_cgroup_sockets_enabled && sk->sk_cgrp)
+ prot = sk->sk_cgrp->sysctl_mem;
    return prot[index];
}

+static inline void memcg_memory_allocated_add(struct cg_proto *prot,
+        unsigned long amt,
+        int *parent_status)
+{
+ struct res_counter *fail;
+ int ret;
+
+ ret = res_counter_charge(prot->memory_allocated,
+ amt << PAGE_SHIFT, &fail);

```

```

+
+ if (ret < 0)
+ *parent_status = OVER_LIMIT;
+}
+
+static inline void memcg_memory_allocated_sub(struct cg_proto *prot,
+      unsigned long amt)
+{
+ res_counter_uncharge(prot->memory_allocated, amt << PAGE_SHIFT);
+}
+
+static inline u64 memcg_memory_allocated_read(struct cg_proto *prot)
+{
+ u64 ret;
+ ret = res_counter_read_u64(prot->memory_allocated, RES_USAGE);
+ return ret >> PAGE_SHIFT;
+}
+
+static inline long
+sk_memory_allocated(const struct sock *sk)
+{
+ struct proto *prot = sk->sk_prot;
+ if (mem_cgroup_sockets_enabled && sk->sk_cgrp)
+ return memcg_memory_allocated_read(sk->sk_cgrp);
+
+ return atomic_long_read(prot->memory_allocated);
+}

static inline long
-sk_memory_allocated_add(struct sock *sk, int amt)
+sk_memory_allocated_add(struct sock *sk, int amt, int *parent_status)
+{
+ struct proto *prot = sk->sk_prot;
+
+ if (mem_cgroup_sockets_enabled && sk->sk_cgrp) {
+ memcg_memory_allocated_add(sk->sk_cgrp, amt, parent_status);
+ /* update the root cgroup regardless */
+ atomic_long_add_return(amt, prot->memory_allocated);
+ return memcg_memory_allocated_read(sk->sk_cgrp);
+ }
+
+ return atomic_long_add_return(amt, prot->memory_allocated);
+}

static inline void
-sk_memory_allocated_sub(struct sock *sk, int amt)
+sk_memory_allocated_sub(struct sock *sk, int amt, int parent_status)
+{

```

```

    struct proto *prot = sk->sk_prot;
+
+ if (mem_cgroup_sockets_enabled && sk->sk_cgrp &&
+     parent_status != OVER_LIMIT) /* Otherwise was uncharged already */
+     memcg_memory_allocated_sub(sk->sk_cgrp, amt);
+
    atomic_long_sub(amt, prot->memory_allocated);
}

static inline void sk_sockets_allocated_dec(struct sock *sk)
{
    struct proto *prot = sk->sk_prot;
+
+ if (mem_cgroup_sockets_enabled && sk->sk_cgrp) {
+     struct cg_proto *cg_proto = sk->sk_cgrp;
+
+     for (; cg_proto; cg_proto = parent_cg_proto(prot, cg_proto))
+         percpu_counter_dec(cg_proto->sockets_allocated);
+ }
+
    percpu_counter_dec(prot->sockets_allocated);
}

static inline void sk_sockets_allocated_inc(struct sock *sk)
{
    struct proto *prot = sk->sk_prot;
+
+ if (mem_cgroup_sockets_enabled && sk->sk_cgrp) {
+     struct cg_proto *cg_proto = sk->sk_cgrp;
+
+     for (; cg_proto; cg_proto = parent_cg_proto(prot, cg_proto))
+         percpu_counter_inc(cg_proto->sockets_allocated);
+ }
+
    percpu_counter_inc(prot->sockets_allocated);
}

@@ -934,6 +1070,9 @@ sk_sockets_allocated_read_positive(struct sock *sk)
{
    struct proto *prot = sk->sk_prot;

+ if (mem_cgroup_sockets_enabled && sk->sk_cgrp)
+     return percpu_counter_sum_positive(sk->sk_cgrp->sockets_allocated);
+
    return percpu_counter_sum_positive(prot->sockets_allocated);
}

diff --git a/mm/memcontrol.c b/mm/memcontrol.c

```



index 3becb24..beedff3 100644

--- a/mm/memcontrol.c

+++ b/mm/memcontrol.c

@@ -379,7 +379,48 @@ enum mem\_type {

```
static void mem_cgroup_get(struct mem_cgroup *memcg);
static void mem_cgroup_put(struct mem_cgroup *memcg);
-static struct mem_cgroup *parent_mem_cgroup(struct mem_cgroup *memcg);
+
+/* Writing them here to avoid exposing memcg's inner layout */
#ifdef CONFIG_CGROUP_MEM_RES_CTLR_KMEM
#ifdef CONFIG_INET
#include <net/sock.h>
+
+static bool mem_cgroup_is_root(struct mem_cgroup *memcg);
+void sock_update_memcg(struct sock *sk)
+{
+ /* A socket spends its whole life in the same cgroup */
+ if (sk->sk_cgrp) {
+ WARN_ON(1);
+ return;
+ }
+ if (static_branch(&memcg_socket_limit_enabled)) {
+ struct mem_cgroup *memcg;
+
+ BUG_ON(!sk->sk_prot->proto_cgroup);
+
+ rcu_read_lock();
+ memcg = mem_cgroup_from_task(current);
+ if (!mem_cgroup_is_root(memcg)) {
+ mem_cgroup_get(memcg);
+ sk->sk_cgrp = sk->sk_prot->proto_cgroup(memcg);
+ }
+ rcu_read_unlock();
+ }
+}
+EXPORT_SYMBOL(sock_update_memcg);
+
+void sock_release_memcg(struct sock *sk)
+{
+ if (static_branch(&memcg_socket_limit_enabled) && sk->sk_cgrp) {
+ struct mem_cgroup *memcg;
+ WARN_ON(!sk->sk_cgrp->memcg);
+ memcg = sk->sk_cgrp->memcg;
+ mem_cgroup_put(memcg);
+ }
+}
+#endif /* CONFIG_INET */
```

```

+ #endif /* CONFIG_CGROUP_MEM_RES_CTLR_KMEM */
+
static void drain_all_stock_async(struct mem_cgroup *memcg);

static struct mem_cgroup_per_zone *
@@ -4930,12 +4971,13 @@ static void mem_cgroup_put(struct mem_cgroup *memcg)
/*
 * Returns the parent mem_cgroup in memcgroup hierarchy with hierarchy enabled.
 */
-static struct mem_cgroup *parent_mem_cgroup(struct mem_cgroup *memcg)
+struct mem_cgroup *parent_mem_cgroup(struct mem_cgroup *memcg)
{
    if (!memcg->res.parent)
        return NULL;
    return mem_cgroup_from_res_counter(memcg->res.parent, res);
}
+EXPORT_SYMBOL(parent_mem_cgroup);

#ifdef CONFIG_CGROUP_MEM_RES_CTLR_SWAP
static void __init enable_swap_cgroup(void)
diff --git a/net/core/sock.c b/net/core/sock.c
index 2b86d24..39e5d01 100644
--- a/net/core/sock.c
+++ b/net/core/sock.c
@@ -111,6 +111,7 @@
#include <linux/init.h>
#include <linux/highmem.h>
#include <linux/user_namespace.h>
+#include <linux/jump_label.h>

#include <asm/uaccess.h>
#include <asm/system.h>
@@ -141,6 +142,9 @@
static struct lock_class_key af_family_keys[AF_MAX];
static struct lock_class_key af_family_slock_keys[AF_MAX];

+struct jump_label_key memcg_socket_limit_enabled;
+EXPORT_SYMBOL(memcg_socket_limit_enabled);
+
/*
 * Make lock validator output more readable. (we pre-construct these
 * strings build-time, so that runtime initialization of socket
@@ -1677,21 +1681,25 @@ int __sk_mem_schedule(struct sock *sk, int size, int kind)
    struct proto *prot = sk->sk_prot;
    int amt = sk_mem_pages(size);
    long allocated;
+ int parent_status = UNDER_LIMIT;

```

```

sk->sk_forward_alloc += amt * SK_MEM_QUANTUM;

- allocated = sk_memory_allocated_add(sk, amt);
+ allocated = sk_memory_allocated_add(sk, amt, &parent_status);

/* Under limit. */
- if (allocated <= sk_prot_mem_limits(sk, 0))
+ if (parent_status == UNDER_LIMIT &&
+ allocated <= sk_prot_mem_limits(sk, 0))
    sk_leave_memory_pressure(sk);

- /* Under pressure. */
- if (allocated > sk_prot_mem_limits(sk, 1))
+ /* Under pressure. (we or our parents) */
+ if ((parent_status > SOFT_LIMIT) ||
+ allocated > sk_prot_mem_limits(sk, 1))
    sk_enter_memory_pressure(sk);

- /* Over hard limit. */
- if (allocated > sk_prot_mem_limits(sk, 2))
+ /* Over hard limit (we or our parents) */
+ if ((parent_status == OVER_LIMIT) ||
+ (allocated > sk_prot_mem_limits(sk, 2)))
    goto suppress_allocation;

/* guarantee minimum buffer size under pressure */
@@ -1738,7 +1746,7 @@ suppress_allocation:
/* Alas. Undo changes. */
sk->sk_forward_alloc -= amt * SK_MEM_QUANTUM;

- sk_memory_allocated_sub(sk, amt);
+ sk_memory_allocated_sub(sk, amt, parent_status);

return 0;
}
@@ -1751,7 +1759,7 @@ EXPORT_SYMBOL(__sk_mem_schedule);
void __sk_mem_reclaim(struct sock *sk)
{
    sk_memory_allocated_sub(sk,
-   sk->sk_forward_alloc >> SK_MEM_QUANTUM_SHIFT);
+   sk->sk_forward_alloc >> SK_MEM_QUANTUM_SHIFT, 0);
    sk->sk_forward_alloc &= SK_MEM_QUANTUM - 1;

    if (sk_under_memory_pressure(sk) &&
--
1.7.6.4

```

---



---

Subject: [PATCH v8 4/9] tcp memory pressure controls  
Posted by [Glauber Costa](#) on Mon, 05 Dec 2011 21:34:58 GMT  
[View Forum Message](#) <> [Reply to Message](#)

---

This patch introduces memory pressure controls for the tcp protocol. It uses the generic socket memory pressure code introduced in earlier patches, and fills in the necessary data in cg\_proto struct.

Signed-off-by: Glauber Costa <glommer@parallels.com>  
CC: KAMEZAWA Hiroyuki <kamezawa.hiroyu@jp.fujitsu.com>  
CC: Eric W. Biederman <ebiederm@xmission.com>

---

```
Documentation/cgroups/memory.txt | 2 +
include/linux/memcontrol.h       | 1 +
include/net/sock.h                | 2 +
include/net/tcp_memcontrol.h      | 17 ++++++++
mm/memcontrol.c                  | 40 ++++++++
net/core/sock.c                   | 43 ++++++++
net/ipv4/Makefile                 | 1 +
net/ipv4/tcp_ipv4.c               | 9 ++++
net/ipv4/tcp_memcontrol.c         | 74 ++++++++
net/ipv6/tcp_ipv6.c               | 5 +++
10 files changed, 189 insertions(+), 5 deletions(-)
create mode 100644 include/net/tcp_memcontrol.h
create mode 100644 net/ipv4/tcp_memcontrol.c
```

```
diff --git a/Documentation/cgroups/memory.txt b/Documentation/cgroups/memory.txt
index 23a8dc5..687dea5 100644
```

```
--- a/Documentation/cgroups/memory.txt
+++ b/Documentation/cgroups/memory.txt
@@ -293,6 +293,8 @@ to trigger slab reclaim when those limits are reached.
thresholds. The Memory Controller allows them to be controlled individually
per cgroup, instead of globally.
```

```
+* tcp memory pressure: sockets memory pressure for the tcp protocol.
```

```
+
```

### 3. User Interface

#### 0. Configuration

```
diff --git a/include/linux/memcontrol.h b/include/linux/memcontrol.h
index f15021b..1513994 100644
```

```
--- a/include/linux/memcontrol.h
+++ b/include/linux/memcontrol.h
@@ -86,6 +86,7 @@ extern struct mem_cgroup *mem_cgroup_from_task(struct task_struct *p);
extern struct mem_cgroup *try_get_mem_cgroup_from_mm(struct mm_struct *mm);

extern struct mem_cgroup *parent_mem_cgroup(struct mem_cgroup *memcg);
+extern struct mem_cgroup *mem_cgroup_from_cont(struct cgroup *cont);
```

```

static inline
int mm_match_cgroup(const struct mm_struct *mm, const struct mem_cgroup *cgroup)
diff --git a/include/net/sock.h b/include/net/sock.h
index b8a63f8..910cb0b 100644
--- a/include/net/sock.h
+++ b/include/net/sock.h
@@ -64,6 +64,8 @@
#include <net/dst.h>
#include <net/checksum.h>

+int mem_cgroup_sockets_init(struct cgroup *cgrp, struct cgroup_subsys *ss);
+void mem_cgroup_sockets_destroy(struct cgroup *cgrp, struct cgroup_subsys *ss);
/*
 * This structure really needs to be cleaned up.
 * Most of it is for TCP, and not used by any of
diff --git a/include/net/tcp_memcontrol.h b/include/net/tcp_memcontrol.h
new file mode 100644
index 0000000..5f5e158
--- /dev/null
+++ b/include/net/tcp_memcontrol.h
@@ -0,0 +1,17 @@
+#ifndef _TCP_MEMCG_H
+#define _TCP_MEMCG_H
+
+
+struct tcp_memcontrol {
+ struct cg_proto cg_proto;
+ /* per-cgroup tcp memory pressure knobs */
+ struct res_counter tcp_memory_allocated;
+ struct percpu_counter tcp_sockets_allocated;
+ /* those two are read-mostly, leave them at the end */
+ long tcp_prot_mem[3];
+ int tcp_memory_pressure;
+};
+
+struct cg_proto *tcp_proto_cgroup(struct mem_cgroup *memcg);
+int tcp_init_cgroup(struct cgroup *cgrp, struct cgroup_subsys *ss);
+void tcp_destroy_cgroup(struct cgroup *cgrp, struct cgroup_subsys *ss);
+#endif /* _TCP_MEMCG_H */
diff --git a/mm/memcontrol.c b/mm/memcontrol.c
index beedff3..b121127 100644
--- a/mm/memcontrol.c
+++ b/mm/memcontrol.c
@@ -50,6 +50,8 @@
#include <linux/cpu.h>
#include <linux/oom.h>
#include "internal.h"
#include <net/sock.h>

```

```

#include <net/tcp_memcontrol.h>

#include <asm/uaccess.h>

@@ -295,6 +297,10 @@ struct mem_cgroup {
    */
    struct mem_cgroup_stat_cpu nocpu_base;
    spinlock_t pcp_counter_lock;
+
+
#ifdef CONFIG_INET
+ struct tcp_memcontrol tcp_mem;
+
#endif
};

/* Stuffs for move charges at task migration. */
@@ -384,6 +390,7 @@ static void mem_cgroup_put(struct mem_cgroup *memcg);
#ifdef CONFIG_CGROUP_MEM_RES_CTLR_KMEM
#ifdef CONFIG_INET
#include <net/sock.h>
#include <net/ip.h>

static bool mem_cgroup_is_root(struct mem_cgroup *memcg);
void sock_update_memcg(struct sock *sk)
@@ -418,6 +425,15 @@ void sock_release_memcg(struct sock *sk)
    mem_cgroup_put(memcg);
}
}
+
+
+struct cg_proto *tcp_proto_cgroup(struct mem_cgroup *memcg)
+{
+ if (!memcg || mem_cgroup_is_root(memcg))
+ return NULL;
+
+ return &memcg->tcp_mem.cg_proto;
+}
+EXPORT_SYMBOL(tcp_proto_cgroup);
#endif /* CONFIG_INET */
#endif /* CONFIG_CGROUP_MEM_RES_CTLR_KMEM */

@@ -800,7 +816,7 @@ static void memcg_check_events(struct mem_cgroup *memcg, struct
page *page)
    preempt_enable();
}

-static struct mem_cgroup *mem_cgroup_from_cont(struct cgroup *cont)
+struct mem_cgroup *mem_cgroup_from_cont(struct cgroup *cont)
{
    return container_of(cgroup_subsys_state(cont,

```

```

    mem_cgroup_subsys_id), struct mem_cgroup,
@@ -4730,14 +4746,34 @@ static int register_kmem_files(struct cgroup *cont, struct
cgroup_subsys *ss)

```

```

    ret = cgroup_add_files(cont, ss, kmem_cgroup_files,
        ARRAY_SIZE(kmem_cgroup_files));
+
+ /*
+  * Part of this would be better living in a separate allocation
+  * function, leaving us with just the cgroup tree population work.
+  * We, however, depend on state such as network's proto_list that
+  * is only initialized after cgroup creation. I found the less
+  * cumbersome way to deal with it to defer it all to populate time
+  */
+ if (!ret)
+   ret = mem_cgroup_sockets_init(cont, ss);
    return ret;
};

```

```

+static void kmem_cgroup_destroy(struct cgroup_subsys *ss,
+ struct cgroup *cont)
+{
+ mem_cgroup_sockets_destroy(cont, ss);
+}
+else
+static int register_kmem_files(struct cgroup *cont, struct cgroup_subsys *ss)
+{
+   return 0;
+}
+
+static void kmem_cgroup_destroy(struct cgroup_subsys *ss,
+ struct cgroup *cont)
+{
+}
+endif

```

```

static struct cftype mem_cgroup_files[] = {
@@ -5096,6 +5132,8 @@ static void mem_cgroup_destroy(struct cgroup_subsys *ss,
{
    struct mem_cgroup *memcg = mem_cgroup_from_cont(cont);

+ kmem_cgroup_destroy(ss, cont);
+
    mem_cgroup_put(memcg);
}

```

```

diff --git a/net/core/sock.c b/net/core/sock.c
index 39e5d01..3d6e370 100644

```

```

--- a/net/core/sock.c
+++ b/net/core/sock.c
@@ -135,6 +135,46 @@
#include <net/tcp.h>
#endif

+static DEFINE_RWLOCK(proto_list_lock);
+static LIST_HEAD(proto_list);
+
+#ifdef CONFIG_CGROUP_MEM_RES_CTLR_KMEM
+int mem_cgroup_sockets_init(struct cgroup *cgrp, struct cgroup_subsys *ss)
+{
+ struct proto *proto;
+ int ret = 0;
+
+ read_lock(&proto_list_lock);
+ list_for_each_entry(proto, &proto_list, node) {
+ if (proto->init_cgroup) {
+ ret = proto->init_cgroup(cgrp, ss);
+ if (ret)
+ goto out;
+ }
+ }
+ read_unlock(&proto_list_lock);
+ return ret;
+out:
+ list_for_each_entry_continue_reverse(proto, &proto_list, node)
+ if (proto->destroy_cgroup)
+ proto->destroy_cgroup(cgrp, ss);
+ read_unlock(&proto_list_lock);
+ return ret;
+}
+
+void mem_cgroup_sockets_destroy(struct cgroup *cgrp, struct cgroup_subsys *ss)
+{
+ struct proto *proto;
+
+ read_lock(&proto_list_lock);
+ list_for_each_entry_reverse(proto, &proto_list, node)
+ if (proto->destroy_cgroup)
+ proto->destroy_cgroup(cgrp, ss);
+ read_unlock(&proto_list_lock);
+}
+#endif
+
+/*
+ * Each address family might have different locking rules, so we have

```



```

* one sock key per address family:
@@ -2256,9 +2296,6 @@ void sk_common_release(struct sock *sk)
}
EXPORT_SYMBOL(sk_common_release);

-static DEFINE_RWLOCK(proto_list_lock);
-static LIST_HEAD(proto_list);
-
#ifdef CONFIG_PROC_FS
#define PROTO_INUSE_NR 64 /* should be enough for the first time */
struct prot_inuse {
diff --git a/net/ipv4/Makefile b/net/ipv4/Makefile
index f2dc69c..dc67a99 100644
--- a/net/ipv4/Makefile
+++ b/net/ipv4/Makefile
@@ -47,6 +47,7 @@ obj-$(CONFIG_TCP_CONG_SCALABLE) += tcp_scalable.o
obj-$(CONFIG_TCP_CONG_LP) += tcp_lp.o
obj-$(CONFIG_TCP_CONG_YEAH) += tcp_yeah.o
obj-$(CONFIG_TCP_CONG_ILLINOIS) += tcp_illinois.o
+obj-$(CONFIG_CGROUP_MEM_RES_CTLR_KMEM) += tcp_memcontrol.o
obj-$(CONFIG_NETLABEL) += cipso_ipv4.o

obj-$(CONFIG_XFRM) += xfrm4_policy.o xfrm4_state.o xfrm4_input.o \
diff --git a/net/ipv4/tcp_ipv4.c b/net/ipv4/tcp_ipv4.c
index d1f4bf8..f70923e 100644
--- a/net/ipv4/tcp_ipv4.c
+++ b/net/ipv4/tcp_ipv4.c
@@ -73,6 +73,7 @@
#include <net/xfrm.h>
#include <net/netdma.h>
#include <net/secure_seq.h>
+#include <net/tcp_memcontrol.h>

#include <linux/inet.h>
#include <linux/ipv6.h>
@@ -1915,6 +1916,7 @@ static int tcp_v4_init_sock(struct sock *sk)
sk->sk_rcvbuf = sysctl_tcp_rmem[1];

local_bh_disable();
+ sock_update_memcg(sk);
sk_sockets_allocated_inc(sk);
local_bh_enable();

@@ -1972,6 +1974,7 @@ void tcp_v4_destroy_sock(struct sock *sk)
}

sk_sockets_allocated_dec(sk);
+ sock_release_memcg(sk);

```

```

}
EXPORT_SYMBOL(tcp_v4_destroy_sock);

@@ -2632,10 +2635,14 @@ struct proto tcp_prot = {
    .compat_setsockopt = compat_tcp_setsockopt,
    .compat_getsockopt = compat_tcp_getsockopt,
    #endif
    #ifdef CONFIG_CGROUP_MEM_RES_CTLR_KMEM
    + .init_cgroup = tcp_init_cgroup,
    + .destroy_cgroup = tcp_destroy_cgroup,
    + .proto_cgroup = tcp_proto_cgroup,
    #endif
};
EXPORT_SYMBOL(tcp_prot);

-
static int __net_init tcp_sk_init(struct net *net)
{
    return inet_ctl_sock_create(&net->ipv4.tcp_sock,
diff --git a/net/ipv4/tcp_memcontrol.c b/net/ipv4/tcp_memcontrol.c
new file mode 100644
index 0000000..4a68d2c
--- /dev/null
+++ b/net/ipv4/tcp_memcontrol.c
@@ -0,0 +1,74 @@
#include <net/tcp.h>
#include <net/tcp_memcontrol.h>
#include <net/sock.h>
#include <linux/memcontrol.h>
#include <linux/module.h>
+
+static inline struct tcp_memcontrol *tcp_from_cgproto(struct cg_proto *cg_proto)
+{
+ return container_of(cg_proto, struct tcp_memcontrol, cg_proto);
+}
+
+static void memcg_tcp_enter_memory_pressure(struct sock *sk)
+{
+ if (!sk->sk_cgrp->memory_pressure)
+ *sk->sk_cgrp->memory_pressure = 1;
+}
+EXPORT_SYMBOL(memcg_tcp_enter_memory_pressure);
+
+int tcp_init_cgroup(struct cgroup *cgrp, struct cgroup_subsys *ss)
+{
+ /*
+ * The root cgroup does not use res_counters, but rather,
+ * rely on the data already collected by the network

```

```

+ * subsystem
+ */
+ struct res_counter *res_parent = NULL;
+ struct cg_proto *cg_proto, *parent_cg;
+ struct tcp_memcontrol *tcp;
+ struct mem_cgroup *memcg = mem_cgroup_from_cont(cgrp);
+ struct mem_cgroup *parent = parent_mem_cgroup(memcg);
+
+ cg_proto = tcp_prot.proto_cgroup(memcg);
+ if (!cg_proto)
+ return 0;
+
+ tcp = tcp_from_cgproto(cg_proto);
+
+ tcp->tcp_prot_mem[0] = sysctl_tcp_mem[0];
+ tcp->tcp_prot_mem[1] = sysctl_tcp_mem[1];
+ tcp->tcp_prot_mem[2] = sysctl_tcp_mem[2];
+ tcp->tcp_memory_pressure = 0;
+
+ parent_cg = tcp_prot.proto_cgroup(parent);
+ if (parent_cg)
+ res_parent = parent_cg->memory_allocated;
+
+ res_counter_init(&tcp->tcp_memory_allocated, res_parent);
+ percpu_counter_init(&tcp->tcp_sockets_allocated, 0);
+
+ cg_proto->enter_memory_pressure = memcg_tcp_enter_memory_pressure;
+ cg_proto->memory_pressure = &tcp->tcp_memory_pressure;
+ cg_proto->sysctl_mem = tcp->tcp_prot_mem;
+ cg_proto->memory_allocated = &tcp->tcp_memory_allocated;
+ cg_proto->sockets_allocated = &tcp->tcp_sockets_allocated;
+ cg_proto->memcg = memcg;
+
+ return 0;
+}
+EXPORT_SYMBOL(tcp_init_cgroup);
+
+void tcp_destroy_cgroup(struct cgroup *cgrp, struct cgroup_subsys *ss)
+{
+ struct mem_cgroup *memcg = mem_cgroup_from_cont(cgrp);
+ struct cg_proto *cg_proto;
+ struct tcp_memcontrol *tcp;
+
+ cg_proto = tcp_prot.proto_cgroup(memcg);
+ if (!cg_proto)
+ return;
+
+ tcp = tcp_from_cgproto(cg_proto);

```

```

+ percpu_counter_destroy(&tcp->tcp_sockets_allocated);
+}
+EXPORT_SYMBOL(tcp_destroy_cgroup);
diff --git a/net/ipv6/tcp_ipv6.c b/net/ipv6/tcp_ipv6.c
index e666768..820ae82 100644
--- a/net/ipv6/tcp_ipv6.c
+++ b/net/ipv6/tcp_ipv6.c
@@ -62,6 +62,7 @@
#include <net/netdma.h>
#include <net/inet_common.h>
#include <net/secure_seq.h>
+#include <net/tcp_memcontrol.h>

#include <asm/uaccess.h>

@@ -1995,6 +1996,7 @@ static int tcp_v6_init_sock(struct sock *sk)
    sk->sk_rcvbuf = sysctl_tcp_rmem[1];

    local_bh_disable();
+ sock_update_memcg(sk);
    sk_sockets_allocated_inc(sk);
    local_bh_enable();

@@ -2228,6 +2230,9 @@ struct proto tcpv6_prot = {
    .compat_setsockopt = compat_tcp_setsockopt,
    .compat_getsockopt = compat_tcp_getsockopt,
#ifdef CONFIG_CGROUP_MEM_RES_CTLR_KMEM
+ .proto_cgroup = tcp_proto_cgroup,
#endif
};

static const struct inet6_protocol tcpv6_protocol = {
--
1.7.6.4

```

---

Subject: [PATCH v8 5/9] per-netns ipv4 sysctl\_tcp\_mem  
 Posted by [Glauber Costa](#) on Mon, 05 Dec 2011 21:34:59 GMT  
[View Forum Message](#) <> [Reply to Message](#)

---

This patch allows each namespace to independently set up its levels for tcp memory pressure thresholds. This patch alone does not buy much: we need to make this values per group of process somehow. This is achieved in the patches that follows in this patchset.

Signed-off-by: Glauber Costa <glommer@parallels.com>

CC: KAMEZAWA Hiroyuki <kamezawa.hiroyu@jp.fujitsu.com>

CC: David S. Miller <davem@davemloft.net>

CC: Eric W. Biederman <ebiederm@xmission.com>

---

```
include/net/netns/ipv4.h | 1 +
include/net/tcp.h        | 1 -
net/ipv4/af_inet.c       | 2 +
net/ipv4/sysctl_net_ipv4.c | 51 ++++++-----
net/ipv4/tcp.c           | 11 +-----
net/ipv4/tcp_ipv4.c      | 1 -
net/ipv4/tcp_memcontrol.c | 9 +++++-
net/ipv6/af_inet6.c      | 2 +
net/ipv6/tcp_ipv6.c      | 1 -
9 files changed, 57 insertions(+), 22 deletions(-)
```

diff --git a/include/net/netns/ipv4.h b/include/net/netns/ipv4.h

index d786b4f..bbd023a 100644

--- a/include/net/netns/ipv4.h

+++ b/include/net/netns/ipv4.h

```
@@ -55,6 +55,7 @@ struct netns_ipv4 {
    int current_rt_cache_rebuild_count;
```

```
    unsigned int sysctl_ping_group_range[2];
+ long sysctl_tcp_mem[3];
```

```
    atomic_t rt_genid;
```

```
    atomic_t dev_addr_genid;
```

diff --git a/include/net/tcp.h b/include/net/tcp.h

index f080e0b..61c7e76 100644

--- a/include/net/tcp.h

+++ b/include/net/tcp.h

```
@@ -230,7 +230,6 @@ extern int sysctl_tcp_fack;
```

```
extern int sysctl_tcp_reordering;
```

```
extern int sysctl_tcp_ecn;
```

```
extern int sysctl_tcp_dsack;
```

```
-extern long sysctl_tcp_mem[3];
```

```
extern int sysctl_tcp_wmem[3];
```

```
extern int sysctl_tcp_rmem[3];
```

```
extern int sysctl_tcp_app_win;
```

diff --git a/net/ipv4/af\_inet.c b/net/ipv4/af\_inet.c

index 1b5096a..a8bbcff 100644

--- a/net/ipv4/af\_inet.c

+++ b/net/ipv4/af\_inet.c

```
@@ -1671,6 +1671,8 @@ static int __init inet_init(void)
```

```
    ip_static_sysctl_init();
```

```
#endif
```

```
+ tcp_prot.sysctl_mem = init_net.ipv4.sysctl_tcp_mem;
```

```

+
+/*
+ * Add all the base protocols.
+ */
diff --git a/net/ipv4/sysctl_net_ipv4.c b/net/ipv4/sysctl_net_ipv4.c
index 69fd720..bbd67ab 100644
--- a/net/ipv4/sysctl_net_ipv4.c
+++ b/net/ipv4/sysctl_net_ipv4.c
@@ -14,6 +14,7 @@
#include <linux/init.h>
#include <linux/slab.h>
#include <linux/nsproxy.h>
+#include <linux/swap.h>
#include <net/snmp.h>
#include <net/icmp.h>
#include <net/ip.h>
@@ -174,6 +175,36 @@ static int proc_allowed_congestion_control(ctl_table *ctl,
return ret;
}

+static int ipv4_tcp_mem(ctl_table *ctl, int write,
+ void __user *buffer, size_t *lenp,
+ loff_t *ppos)
+{
+ int ret;
+ unsigned long vec[3];
+ struct net *net = current->nsproxy->net_ns;
+
+ ctl_table tmp = {
+ .data = &vec,
+ .maxlen = sizeof(vec),
+ .mode = ctl->mode,
+ };
+
+ if (!write) {
+ ctl->data = &net->ipv4.sysctl_tcp_mem;
+ return proc_doulongvec_minmax(ctl, write, buffer, lenp, ppos);
+ }
+
+ ret = proc_doulongvec_minmax(&tmp, write, buffer, lenp, ppos);
+ if (ret)
+ return ret;
+
+ net->ipv4.sysctl_tcp_mem[0] = vec[0];
+ net->ipv4.sysctl_tcp_mem[1] = vec[1];
+ net->ipv4.sysctl_tcp_mem[2] = vec[2];
+
+ return 0;

```

```

+}
+
static struct ctl_table ipv4_table[] = {
{
    .procname = "tcp_timestamps",
@@ -433,13 +464,6 @@ static struct ctl_table ipv4_table[] = {
    .proc_handler = proc_dointvec
},
{
- .procname = "tcp_mem",
- .data = &sysctl_tcp_mem,
- .maxlen = sizeof(sysctl_tcp_mem),
- .mode = 0644,
- .proc_handler = proc_doulongvec_minmax
- },
- {
    .procname = "tcp_wmem",
    .data = &sysctl_tcp_wmem,
    .maxlen = sizeof(sysctl_tcp_wmem),
@@ -721,6 +745,12 @@ static struct ctl_table ipv4_net_table[] = {
    .mode = 0644,
    .proc_handler = ipv4_ping_group_range,
},
+ {
+ .procname = "tcp_mem",
+ .maxlen = sizeof(init_net.ipv4.sysctl_tcp_mem),
+ .mode = 0644,
+ .proc_handler = ipv4_tcp_mem,
+ },
+ { }
};

@@ -734,6 +764,7 @@ EXPORT_SYMBOL_GPL(net_ipv4_ctl_path);
static __net_init int ipv4_sysctl_init_net(struct net *net)
{
    struct ctl_table *table;
+ unsigned long limit;

    table = ipv4_net_table;
    if (!net_eq(net, &init_net)) {
@@ -769,6 +800,12 @@ static __net_init int ipv4_sysctl_init_net(struct net *net)

    net->ipv4.sysctl_rt_cache_rebuild_count = 4;

+ limit = nr_free_buffer_pages() / 8;
+ limit = max(limit, 128UL);
+ net->ipv4.sysctl_tcp_mem[0] = limit / 4 * 3;
+ net->ipv4.sysctl_tcp_mem[1] = limit;

```

```

+ net->ipv4.sysctl_tcp_mem[2] = net->ipv4.sysctl_tcp_mem[0] * 2;
+
+ net->ipv4.ipv4_hdr = register_net_sysctl_table(net,
+ net_ipv4_ctl_path, table);
+ if (net->ipv4.ipv4_hdr == NULL)
diff --git a/net/ipv4/tcp.c b/net/ipv4/tcp.c
index 34f5db1..5f618d1 100644
--- a/net/ipv4/tcp.c
+++ b/net/ipv4/tcp.c
@@ -282,11 +282,9 @@ int sysctl_tcp_fin_timeout __read_mostly = TCP_FIN_TIMEOUT;
struct percpu_counter tcp_orphan_count;
EXPORT_SYMBOL_GPL(tcp_orphan_count);

-long sysctl_tcp_mem[3] __read_mostly;
-int sysctl_tcp_wmem[3] __read_mostly;
-int sysctl_tcp_rmem[3] __read_mostly;

-EXPORT_SYMBOL(sysctl_tcp_mem);
EXPORT_SYMBOL(sysctl_tcp_rmem);
EXPORT_SYMBOL(sysctl_tcp_wmem);

@@ -3272,14 +3270,9 @@ void __init tcp_init(void)
sysctl_tcp_max_orphans = cnt / 2;
sysctl_max_syn_backlog = max(128, cnt / 256);

- limit = nr_free_buffer_pages() / 8;
- limit = max(limit, 128UL);
- sysctl_tcp_mem[0] = limit / 4 * 3;
- sysctl_tcp_mem[1] = limit;
- sysctl_tcp_mem[2] = sysctl_tcp_mem[0] * 2;
-
/* Set per-socket limits to no more than 1/128 the pressure threshold */
- limit = ((unsigned long)sysctl_tcp_mem[1]) << (PAGE_SHIFT - 7);
+ limit = ((unsigned long)init_net.ipv4.sysctl_tcp_mem[1])
+ << (PAGE_SHIFT - 7);
max_share = min(4UL*1024*1024, limit);

sysctl_tcp_wmem[0] = SK_MEM_QUANTUM;
diff --git a/net/ipv4/tcp_ipv4.c b/net/ipv4/tcp_ipv4.c
index f70923e..cbba5aa 100644
--- a/net/ipv4/tcp_ipv4.c
+++ b/net/ipv4/tcp_ipv4.c
@@ -2621,7 +2621,6 @@ struct proto tcp_prot = {
.orphan_count = &tcp_orphan_count,
.memory_allocated = &tcp_memory_allocated,
.memory_pressure = &tcp_memory_pressure,
- .sysctl_mem = sysctl_tcp_mem,
- .sysctl_wmem = sysctl_tcp_wmem,

```



```

.sysctl_rmem = sysctl_tcp_rmem,
.max_header = MAX_TCP_HEADER,
diff --git a/net/ipv4/tcp_memcontrol.c b/net/ipv4/tcp_memcontrol.c
index 4a68d2c..bfb0c2b 100644
--- a/net/ipv4/tcp_memcontrol.c
+++ b/net/ipv4/tcp_memcontrol.c
@@ -1,6 +1,8 @@
#include <net/tcp.h>
#include <net/tcp_memcontrol.h>
#include <net/sock.h>
+#include <net/ip.h>
+#include <linux/nsproxy.h>
#include <linux/memcontrol.h>
#include <linux/module.h>

@@ -28,6 +30,7 @@ int tcp_init_cgroup(struct cgroup *cgrp, struct cgroup_subsys *ss)
    struct tcp_memcontrol *tcp;
    struct mem_cgroup *memcg = mem_cgroup_from_cont(cgrp);
    struct mem_cgroup *parent = parent_mem_cgroup(memcg);
+ struct net *net = current->nsproxy->net_ns;

    cg_proto = tcp_prot.proto_cgroup(memcg);
    if (!cg_proto)
@@ -35,9 +38,9 @@ int tcp_init_cgroup(struct cgroup *cgrp, struct cgroup_subsys *ss)

    tcp = tcp_from_cgproto(cg_proto);

- tcp->tcp_prot_mem[0] = sysctl_tcp_mem[0];
- tcp->tcp_prot_mem[1] = sysctl_tcp_mem[1];
- tcp->tcp_prot_mem[2] = sysctl_tcp_mem[2];
+ tcp->tcp_prot_mem[0] = net->ipv4.sysctl_tcp_mem[0];
+ tcp->tcp_prot_mem[1] = net->ipv4.sysctl_tcp_mem[1];
+ tcp->tcp_prot_mem[2] = net->ipv4.sysctl_tcp_mem[2];
    tcp->tcp_memory_pressure = 0;

    parent_cg = tcp_prot.proto_cgroup(parent);
diff --git a/net/ipv6/af_inet6.c b/net/ipv6/af_inet6.c
index d27c797..49b2145 100644
--- a/net/ipv6/af_inet6.c
+++ b/net/ipv6/af_inet6.c
@@ -1115,6 +1115,8 @@ static int __init inet6_init(void)
    if (err)
        goto static_sysctl_fail;
    #endif
+ tcpv6_prot.sysctl_mem = init_net.ipv4.sysctl_tcp_mem;
+
+ /*
+  * ipngwg API draft makes clear that the correct semantics

```

\* for TCP and UDP is to consider one TCP and UDP instance

diff --git a/net/ipv6/tcp\_ipv6.c b/net/ipv6/tcp\_ipv6.c

index 820ae82..51bbfb0 100644

--- a/net/ipv6/tcp\_ipv6.c

+++ b/net/ipv6/tcp\_ipv6.c

@@ -2216,7 +2216,6 @@ struct proto tcpv6\_prot = {

.memory\_allocated = &tcp\_memory\_allocated,

.memory\_pressure = &tcp\_memory\_pressure,

.orphan\_count = &tcp\_orphan\_count,

- .sysctl\_mem = sysctl\_tcp\_mem,

.sysctl\_wmem = sysctl\_tcp\_wmem,

.sysctl\_rmem = sysctl\_tcp\_rmem,

.max\_header = MAX\_TCP\_HEADER,

--

1.7.6.4

---

Subject: [PATCH v8 6/9] tcp buffer limitation: per-cgroup limit

Posted by [Glauber Costa](#) on Mon, 05 Dec 2011 21:35:00 GMT

[View Forum Message](#) <> [Reply to Message](#)

---

This patch uses the "tcp.limit\_in\_bytes" field of the kmem\_cgroup to effectively control the amount of kernel memory pinned by a cgroup.

This value is ignored in the root cgroup, and in all others, caps the value specified by the admin in the net namespaces' view of tcp\_sysctl\_mem.

If namespaces are being used, the admin is allowed to set a value bigger than cgroup's maximum, the same way it is allowed to set pretty much unlimited values in a real box.

Signed-off-by: Glauber Costa <glommer@parallels.com>

CC: David S. Miller <davem@davemloft.net>

CC: Hiroyouki Kamezawa <kamezawa.hiroyu@jp.fujitsu.com>

CC: Eric W. Biederman <ebiederm@xmission.com>

---

Documentation/cgroups/memory.txt | 1 +

include/net/tcp\_memcontrol.h | 2 +

net/ipv4/sysctl\_net\_ipv4.c | 14 +++++

net/ipv4/tcp\_memcontrol.c | 137 ++++++

4 files changed, 152 insertions(+), 2 deletions(-)

diff --git a/Documentation/cgroups/memory.txt b/Documentation/cgroups/memory.txt

index 687dea5..1c9779a 100644

--- a/Documentation/cgroups/memory.txt

+++ b/Documentation/cgroups/memory.txt

@@ -78,6 +78,7 @@ Brief summary of control files.



```

+ tcp_prot_mem(memcg, vec[0], 0);
+ tcp_prot_mem(memcg, vec[1], 1);
+ tcp_prot_mem(memcg, vec[2], 2);
+ rcu_read_unlock();
+#endif
+
+ net->ipv4.sysctl_tcp_mem[0] = vec[0];
+ net->ipv4.sysctl_tcp_mem[1] = vec[1];
+ net->ipv4.sysctl_tcp_mem[2] = vec[2];
diff --git a/net/ipv4/tcp_memcontrol.c b/net/ipv4/tcp_memcontrol.c
index bfb0c2b..e353390 100644
--- a/net/ipv4/tcp_memcontrol.c
+++ b/net/ipv4/tcp_memcontrol.c
@@ -6,6 +6,19 @@
#include <linux/memcontrol.h>
#include <linux/module.h>

+static u64 tcp_cgroup_read(struct cgroup *cont, struct cftype *cft);
+static int tcp_cgroup_write(struct cgroup *cont, struct cftype *cft,
+    const char *buffer);
+
+static struct cftype tcp_files[] = {
+ {
+ .name = "kmem.tcp.limit_in_bytes",
+ .write_string = tcp_cgroup_write,
+ .read_u64 = tcp_cgroup_read,
+ .private = RES_LIMIT,
+ },
+};
+
+static inline struct tcp_memcontrol *tcp_from_cgproto(struct cg_proto *cg_proto)
+ {
+     return container_of(cg_proto, struct tcp_memcontrol, cg_proto);
@@ -34,7 +47,7 @@ int tcp_init_cgroup(struct cgroup *cgrp, struct cgroup_subsys *ss)

cg_proto = tcp_prot.proto_cgroup(memcg);
if (!cg_proto)
- return 0;
+ goto create_files;

tcp = tcp_from_cgproto(cg_proto);

@@ -57,7 +70,9 @@ int tcp_init_cgroup(struct cgroup *cgrp, struct cgroup_subsys *ss)
cg_proto->sockets_allocated = &tcp->tcp_sockets_allocated;
cg_proto->memcg = memcg;

- return 0;
+create_files:

```

```

+ return cgroup_add_files(cgrp, ss, tcp_files,
+   ARRAY_SIZE(tcp_files));
}
EXPORT_SYMBOL(tcp_init_cgroup);

@@ -66,6 +81,7 @@ void tcp_destroy_cgroup(struct cgroup *cgrp, struct cgroup_subsys *ss)
    struct mem_cgroup *memcg = mem_cgroup_from_cont(cgrp);
    struct cg_proto *cg_proto;
    struct tcp_memcontrol *tcp;
+ u64 val;

    cg_proto = tcp_prot.proto_cgroup(memcg);
    if (!cg_proto)
@@ -73,5 +89,122 @@ void tcp_destroy_cgroup(struct cgroup *cgrp, struct cgroup_subsys *ss)

    tcp = tcp_from_cgproto(cg_proto);
    percpu_counter_destroy(&tcp->tcp_sockets_allocated);
+
+ val = res_counter_read_u64(&tcp->tcp_memory_allocated, RES_USAGE);
+
+ if (val != RESOURCE_MAX)
+   jump_label_dec(&memcg_socket_limit_enabled);
+
EXPORT_SYMBOL(tcp_destroy_cgroup);
+
+static int tcp_update_limit(struct mem_cgroup *memcg, u64 val)
+{
+ struct net *net = current->nsproxy->net_ns;
+ struct tcp_memcontrol *tcp;
+ struct cg_proto *cg_proto;
+ u64 old_lim;
+ int i;
+ int ret;
+
+ cg_proto = tcp_prot.proto_cgroup(memcg);
+ if (!cg_proto)
+   return -EINVAL;
+
+ if (val > RESOURCE_MAX)
+   val = RESOURCE_MAX;
+
+ tcp = tcp_from_cgproto(cg_proto);
+
+ old_lim = res_counter_read_u64(&tcp->tcp_memory_allocated, RES_LIMIT);
+ ret = res_counter_set_limit(&tcp->tcp_memory_allocated, val);
+ if (ret)
+   return ret;
+

```

```

+ for (i = 0; i < 3; i++)
+ tcp->tcp_prot_mem[i] = min_t(long, val >> PAGE_SHIFT,
+     net->ipv4.sysctl_tcp_mem[i]);
+
+ if (val == RESOURCE_MAX && old_lim != RESOURCE_MAX)
+     jump_label_dec(&memcg_socket_limit_enabled);
+ else if (old_lim == RESOURCE_MAX && val != RESOURCE_MAX)
+     jump_label_inc(&memcg_socket_limit_enabled);
+
+ return 0;
+}
+
+static int tcp_cgroup_write(struct cgroup *cont, struct cftype *cft,
+    const char *buffer)
+{
+ struct mem_cgroup *memcg = mem_cgroup_from_cont(cont);
+ unsigned long long val;
+ int ret = 0;
+
+ switch (cft->private) {
+ case RES_LIMIT:
+     /* see memcontrol.c */
+     ret = res_counter_memparse_write_strategy(buffer, &val);
+     if (ret)
+         break;
+     ret = tcp_update_limit(memcg, val);
+     break;
+ default:
+     ret = -EINVAL;
+     break;
+ }
+ return ret;
+}
+
+static u64 tcp_read_stat(struct mem_cgroup *memcg, int type, u64 default_val)
+{
+ struct tcp_memcontrol *tcp;
+ struct cg_proto *cg_proto;
+
+ cg_proto = tcp_prot.proto_cgroup(memcg);
+ if (!cg_proto)
+     return default_val;
+
+ tcp = tcp_from_cgproto(cg_proto);
+ return res_counter_read_u64(&tcp->tcp_memory_allocated, type);
+}
+
+static u64 tcp_cgroup_read(struct cgroup *cont, struct cftype *cft)

```

```

+{
+ struct mem_cgroup *memcg = mem_cgroup_from_cont(cont);
+ u64 val;
+
+ switch (cft->private) {
+ case RES_LIMIT:
+ val = tcp_read_stat(memcg, RES_LIMIT, RESOURCE_MAX);
+ break;
+ default:
+ BUG();
+ }
+ return val;
+}
+
+unsigned long long tcp_max_memory(const struct mem_cgroup *memcg)
+{
+ struct tcp_memcontrol *tcp;
+ struct cg_proto *cg_proto;
+
+ cg_proto = tcp_prot.proto_cgroup((struct mem_cgroup *)memcg);
+ if (!cg_proto)
+ return 0;
+
+ tcp = tcp_from_cgproto(cg_proto);
+ return res_counter_read_u64(&tcp->tcp_memory_allocated, RES_LIMIT);
+}
+
+void tcp_prot_mem(struct mem_cgroup *memcg, long val, int idx)
+{
+ struct tcp_memcontrol *tcp;
+ struct cg_proto *cg_proto;
+
+ cg_proto = tcp_prot.proto_cgroup(memcg);
+ if (!cg_proto)
+ return;
+
+ tcp = tcp_from_cgproto(cg_proto);
+
+ tcp->tcp_prot_mem[idx] = val;
+}
--
1.7.6.4

```

---

Subject: [PATCH v8 7/9] Display current tcp memory allocation in kmem cgroup  
Posted by [Glauber Costa](#) on Mon, 05 Dec 2011 21:35:01 GMT  
[View Forum Message](#) <> [Reply to Message](#)

---

This patch introduces `kmem.tcp.usage_in_bytes` file, living in the `kmem_cgroup` filesystem. It is a simple read-only file that displays the amount of kernel memory currently consumed by the cgroup.

Signed-off-by: Glauber Costa <glommer@parallels.com>

Reviewed-by: Hiroyouki Kamezawa <kamezawa.hiroyu@jp.fujitsu.com>

CC: David S. Miller <davem@davemloft.net>

CC: Eric W. Biederman <ebiederm@xmission.com>

---

```
Documentation/cgroups/memory.txt | 1 +
net/ipv4/tcp_memcontrol.c         | 21 +++++
2 files changed, 22 insertions(+), 0 deletions(-)
```

```
diff --git a/Documentation/cgroups/memory.txt b/Documentation/cgroups/memory.txt
index 1c9779a..6922b6c 100644
```

```
--- a/Documentation/cgroups/memory.txt
```

```
+++ b/Documentation/cgroups/memory.txt
```

```
@@ -79,6 +79,7 @@ Brief summary of control files.
```

```
memory.independent_kmem_limit # select whether or not kernel memory limits are
independent of user limits
```

```
memory.kmem.tcp.limit_in_bytes # set/show hard limit for tcp buf memory
```

```
+ memory.kmem.tcp.usage_in_bytes # show current tcp buf memory allocation
```

## 1. History

```
diff --git a/net/ipv4/tcp_memcontrol.c b/net/ipv4/tcp_memcontrol.c
```

```
index e353390..9481f23 100644
```

```
--- a/net/ipv4/tcp_memcontrol.c
```

```
+++ b/net/ipv4/tcp_memcontrol.c
```

```
@@ -17,6 +17,11 @@ static struct cftype tcp_files[] = {
```

```
    .read_u64 = tcp_cgroup_read,
```

```
    .private = RES_LIMIT,
```

```
},
```

```
+ {
```

```
+ .name = "kmem.tcp.usage_in_bytes",
```

```
+ .read_u64 = tcp_cgroup_read,
```

```
+ .private = RES_USAGE,
```

```
+ },
```

```
};
```

```
static inline struct tcp_memcontrol *tcp_from_cgproto(struct cg_proto *cg_proto)
```

```
@@ -167,6 +172,19 @@ static u64 tcp_read_stat(struct mem_cgroup *memcg, int type, u64
default_val)
```

```
    return res_counter_read_u64(&tcp->tcp_memory_allocated, type);
```

```
}
```

```
+static u64 tcp_read_usage(struct mem_cgroup *memcg)
```

```
+{
```



```

+ struct tcp_memcontrol *tcp;
+ struct cg_proto *cg_proto;
+
+ cg_proto = tcp_prot.proto_cgroup(memcg);
+ if (!cg_proto)
+ return atomic_long_read(&tcp_memory_allocated) << PAGE_SHIFT;
+
+ tcp = tcp_from_cgproto(cg_proto);
+ return res_counter_read_u64(&tcp->tcp_memory_allocated, RES_USAGE);
+}
+
static u64 tcp_cgroup_read(struct cgroup *cont, struct cftype *cft)
{
    struct mem_cgroup *memcg = mem_cgroup_from_cont(cont);
@@ -176,6 +194,9 @@ static u64 tcp_cgroup_read(struct cgroup *cont, struct cftype *cft)
    case RES_LIMIT:
        val = tcp_read_stat(memcg, RES_LIMIT, RESOURCE_MAX);
        break;
+ case RES_USAGE:
+     val = tcp_read_usage(memcg);
+     break;
    default:
        BUG();
}
--
1.7.6.4

```

---

Subject: [PATCH v8 8/9] Display current tcp failcnt in kmem cgroup  
 Posted by [Glauber Costa](#) on Mon, 05 Dec 2011 21:35:02 GMT  
[View Forum Message](#) <> [Reply to Message](#)

---

This patch introduces kmem.tcp.failcnt file, living in the kmem\_cgroup filesystem. Following the pattern in the other memcg resources, this files keeps a counter of how many times allocation failed due to limits being hit in this cgroup. The root cgroup will always show a failcnt of 0.

Signed-off-by: Glauber Costa <glommer@parallels.com>  
 Reviewed-by: Hiroyouki Kamezawa <kamezawa.hiroyu@jp.fujitsu.com>  
 CC: David S. Miller <davem@davemloft.net>  
 CC: Eric W. Biederman <ebiederm@xmission.com>

---

```

net/ipv4/tcp_memcontrol.c | 31 +++++
1 files changed, 31 insertions(+), 0 deletions(-)

```

```

diff --git a/net/ipv4/tcp_memcontrol.c b/net/ipv4/tcp_memcontrol.c
index 9481f23..d438fba 100644

```

```

--- a/net/ipv4/tcp_memcontrol.c
+++ b/net/ipv4/tcp_memcontrol.c
@@ -9,6 +9,7 @@
static u64 tcp_cgroup_read(struct cgroup *cont, struct cftype *cft);
static int tcp_cgroup_write(struct cgroup *cont, struct cftype *cft,
    const char *buffer);
+static int tcp_cgroup_reset(struct cgroup *cont, unsigned int event);

static struct cftype tcp_files[] = {
{
@@ -22,6 +23,12 @@ static struct cftype tcp_files[] = {
    .read_u64 = tcp_cgroup_read,
    .private = RES_USAGE,
},
+ {
+ .name = "kmem.tcp.failcnt",
+ .private = RES_FAILCNT,
+ .trigger = tcp_cgroup_reset,
+ .read_u64 = tcp_cgroup_read,
+ },
};

static inline struct tcp_memcontrol *tcp_from_cgproto(struct cg_proto *cg_proto)
@@ -197,12 +204,36 @@ static u64 tcp_cgroup_read(struct cgroup *cont, struct cftype *cft)
    case RES_USAGE:
        val = tcp_read_usage(memcg);
        break;
+ case RES_FAILCNT:
+     val = tcp_read_stat(memcg, RES_FAILCNT, 0);
+     break;
    default:
        BUG();
}
return val;
}

+static int tcp_cgroup_reset(struct cgroup *cont, unsigned int event)
+{
+ struct mem_cgroup *memcg;
+ struct tcp_memcontrol *tcp;
+ struct cg_proto *cg_proto;
+
+ memcg = mem_cgroup_from_cont(cont);
+ cg_proto = tcp_prot.proto_cgroup(memcg);
+ if (!cg_proto)
+     return 0;
+ tcp = tcp_from_cgproto(cg_proto);
+

```

```

+ switch (event) {
+ case RES_FAILCNT:
+ res_counter_reset_failcnt(&tcp->tcp_memory_allocated);
+ break;
+ }
+
+ return 0;
+}
+
+ unsigned long long tcp_max_memory(const struct mem_cgroup *memcg)
+ {
+ struct tcp_memcontrol *tcp;
+
+ --
+ 1.7.6.4

```

---

Subject: [PATCH v8 9/9] Display maximum tcp memory allocation in kmem cgroup  
 Posted by [Glauber Costa](#) on Mon, 05 Dec 2011 21:35:03 GMT  
[View Forum Message](#) <> [Reply to Message](#)

---

This patch introduces kmem.tcp.max\_usage\_in\_bytes file, living in the kmem\_cgroup filesystem. The root cgroup will display a value equal to RESOURCE\_MAX. This is to avoid introducing any locking schemes in the network paths when cgroups are not being actively used.

All others, will see the maximum memory ever used by this cgroup.

Signed-off-by: Glauber Costa <glommer@parallels.com>  
 Reviewed-by: Hiroyouki Kamezawa <kamezawa.hiroyu@jp.fujitsu.com>  
 CC: David S. Miller <davem@davemloft.net>  
 CC: Eric W. Biederman <ebiederm@xmission.com>

---

```

net/ipv4/tcp_memcontrol.c | 12 ++++++++
1 files changed, 11 insertions(+), 1 deletions(-)

```

```

diff --git a/net/ipv4/tcp_memcontrol.c b/net/ipv4/tcp_memcontrol.c
index d438fba..171d7b6 100644
--- a/net/ipv4/tcp_memcontrol.c
+++ b/net/ipv4/tcp_memcontrol.c
@@ -29,6 +29,12 @@ static struct cftype tcp_files[] = {
     .trigger = tcp_cgroup_reset,
     .read_u64 = tcp_cgroup_read,
 },
+ {
+ .name = "kmem.tcp.max_usage_in_bytes",
+ .private = RES_MAX_USAGE,
+ .trigger = tcp_cgroup_reset,
+ .read_u64 = tcp_cgroup_read,

```

```

+ },
+ };

static inline struct tcp_memcontrol *tcp_from_cgproto(struct cg_proto *cg_proto)
@@ -205,7 +211,8 @@ static u64 tcp_cgroup_read(struct cgroup *cont, struct cftype *cft)
    val = tcp_read_usage(memcg);
    break;
    case RES_FAILCNT:
-   val = tcp_read_stat(memcg, RES_FAILCNT, 0);
+ case RES_MAX_USAGE:
+   val = tcp_read_stat(memcg, cft->private, 0);
    break;
    default:
    BUG();
@@ -226,6 +233,9 @@ static int tcp_cgroup_reset(struct cgroup *cont, unsigned int event)
    tcp = tcp_from_cgproto(cg_proto);

    switch (event) {
+ case RES_MAX_USAGE:
+   res_counter_reset_max(&tcp->tcp_memory_allocated);
+   break;
    case RES_FAILCNT:
    res_counter_reset_failcnt(&tcp->tcp_memory_allocated);
    break;
--
1.7.6.4

```

---

Subject: Re: [PATCH v8 0/9] per-cgroup tcp memory pressure controls  
 Posted by [Glauber Costa](#) on Wed, 07 Dec 2011 11:06:17 GMT  
[View Forum Message](#) <> [Reply to Message](#)

---

On 12/05/2011 07:34 PM, Glauber Costa wrote:

```

> Hi,
>
> This is my new attempt to fix all the concerns that were raised during
> the last iteration.
>
> I should highlight:
> 1) proc information is kept intact. (although I kept the wrapper functions)
>   it will be submitted as a follow up patch so it can get the attention it
>   deserves
> 2) sockets now hold a reference to memcg. sockets can be alive even after the
>   task is gone, so we don't bother with between cgroups movements.
>   To be able to release resources more easily in this scenario, the parent
>   pointer in struct cg_proto was replaced by a memcg object. We then iterate
>   through its pointer (which is cleaner anyway)
>

```

> The rest should be mostly the same except for small fixes and style changes.  
>

Kame,

Does this one address your previous concerns?

Thanks

---

---

Subject: Re: [PATCH v8 0/9] per-cgroup tcp memory pressure controls  
Posted by [KAMEZAWA Hiroyuki](#) on Fri, 09 Dec 2011 01:04:20 GMT

[View Forum Message](#) <> [Reply to Message](#)

---

On Wed, 7 Dec 2011 09:06:17 -0200

Glauber Costa <[glommer@parallels.com](mailto:glommer@parallels.com)> wrote:

> On 12/05/2011 07:34 PM, Glauber Costa wrote:

> > Hi,

> >

> > This is my new attempt to fix all the concerns that were raised during  
> > the last iteration.

> >

> > I should highlight:

> > 1) proc information is kept intact. (although I kept the wrapper functions)

> > it will be submitted as a follow up patch so it can get the attention it

> > deserves

> > 2) sockets now hold a reference to memcg. sockets can be alive even after the  
> > task is gone, so we don't bother with between cgroups movements.

> > To be able to release resources more easily in this cenario, the parent

> > pointer in struct cg\_proto was replaced by a memcg object. We then iterate

> > through its pointer (which is cleaner anyway)

> >

> > The rest should be mostly the same except for small fixes and style changes.

> >

>

> Kame,

>

> Does this one address your previous concerns?

>

Your highlight seems good. I'll look into details.

Thanks,

-Kame

---

---

Subject: Re: [PATCH v8 1/9] Basic kernel memory functionality for the Memory

## Controller

Posted by [KAMEZAWA Hiroyuki](#) on Fri, 09 Dec 2011 01:21:13 GMT

[View Forum Message](#) <> [Reply to Message](#)

---

On Mon, 5 Dec 2011 19:34:55 -0200

Glauber Costa <glommer@parallels.com> wrote:

> This patch lays down the foundation for the kernel memory component  
> of the Memory Controller.

>

> As of today, I am only laying down the following files:

>

> \* memory.independent\_kmem\_limit

> \* memory.kmem.limit\_in\_bytes (currently ignored)

> \* memory.kmem.usage\_in\_bytes (always zero)

>

> Signed-off-by: Glauber Costa <glommer@parallels.com>

> Reviewed-by: Kirill A. Shutemov <kirill@shutemov.name>

> CC: Paul Menage <paul@paulmenage.org>

> CC: Greg Thelen <gthelen@google.com>

As I wrote, please CC Johannes and Michal Hocko for memcg related parts.

A few questions.

==

> + val = !!val;

> +

> + if (parent && parent->use\_hierarchy &&

> + (val != parent->kmem\_independent\_accounting))

> + return -EINVAL;

==

Hm, why you check val != parent->kmem\_independent\_accounting ?

if (parent && parent->use\_hierarchy)

return -EINVAL;

?

BTW, you didn't check this cgroup has children or not.

I think

if (this\_cgroup->use\_hierarchy &&

!list\_empty(this\_cgroup->childlen))

return -EINVAL;

==

> + /\*

> + \* TODO: We need to handle the case in which we are doing

> + \* independent kmem accounting as authorized by our parent,

```
> + * but then our parent changes its parameter.  
> + */  
> + cgroup_lock();  
> + memcg->kmem_independent_accounting = val;  
> + cgroup_unlock();
```

Do we need cgroup\_lock() here ?

Thanks,  
-Kame

---

Subject: Re: [PATCH v8 2/9] foundations of per-cgroup memory pressure controlling.

Posted by [KAMEZAWA Hiroyuki](#) on Fri, 09 Dec 2011 01:24:38 GMT

[View Forum Message](#) <> [Reply to Message](#)

---

On Mon, 5 Dec 2011 19:34:56 -0200

Glauber Costa <glommer@parallels.com> wrote:

```
> This patch replaces all uses of struct sock fields' memory_pressure,  
> memory_allocated, sockets_allocated, and sysctl_mem to accessor  
> macros. Those macros can either receive a socket argument, or a mem_cgroup  
> argument, depending on the context they live in.  
>  
> Since we're only doing a macro wrapping here, no performance impact at all is  
> expected in the case where we don't have cgroups disabled.  
>  
> Signed-off-by: Glauber Costa <glommer@parallels.com>  
> CC: David S. Miller <davem@davemloft.net>  
> CC: Hiroyuki Kamezawa <kamezawa.hiroyu@jp.fujitsu.com>  
> CC: Eric W. Biederman <ebiederm@xmission.com>  
> CC: Eric Dumazet <eric.dumazet@gmail.com>
```

please get ack from network guys.  
from me.

Reviewed-by: KAMEZAWA Hiroyuki <kamezawa.hiroyu@jp.fujitsu.com>

---

Subject: Re: [PATCH v8 3/9] socket: initial cgroup code.

Posted by [KAMEZAWA Hiroyuki](#) on Fri, 09 Dec 2011 01:49:01 GMT

[View Forum Message](#) <> [Reply to Message](#)

---

On Mon, 5 Dec 2011 19:34:57 -0200

Glauber Costa <glommer@parallels.com> wrote:

```
> The goal of this work is to move the memory pressure tcp
```

> controls to a cgroup, instead of just relying on global  
> conditions.  
>  
> To avoid excessive overhead in the network fast paths,  
> the code that accounts allocated memory to a cgroup is  
> hidden inside a static\_branch(). This branch is patched out  
> until the first non-root cgroup is created. So when nobody  
> is using cgroups, even if it is mounted, no significant performance  
> penalty should be seen.  
>  
> This patch handles the generic part of the code, and has nothing  
> tcp-specific.  
>  
> Signed-off-by: Glauber Costa <glommer@parallels.com>  
> CC: Kirill A. Shutemov <kirill@shutemov.name>  
> CC: KAMEZAWA Hiroyuki <kamezawa.hiroyu@jp.fujitsu.com>  
> CC: David S. Miller <davem@davemloft.net>  
> CC: Eric W. Biederman <ebiederm@xmission.com>  
> CC: Eric Dumazet <eric.dumazet@gmail.com>

Reviewed-by: KAMEZAWA Hiroyuki <kamezawa.hiroyu@jp.fujitsu.com>

---

---

Subject: Re: [PATCH v8 4/9] tcp memory pressure controls  
Posted by [KAMEZAWA Hiroyuki](#) on Fri, 09 Dec 2011 01:51:07 GMT  
[View Forum Message](#) <> [Reply to Message](#)

---

On Mon, 5 Dec 2011 19:34:58 -0200  
Glauber Costa <glommer@parallels.com> wrote:

> This patch introduces memory pressure controls for the tcp  
> protocol. It uses the generic socket memory pressure code  
> introduced in earlier patches, and fills in the  
> necessary data in cg\_proto struct.  
>  
> Signed-off-by: Glauber Costa <glommer@parallels.com>  
> CC: KAMEZAWA Hiroyuki <kamezawa.hiroyu@jp.fujitsu.com>  
> CC: Eric W. Biederman <ebiederm@xmission.com>

Reviewed-by: KAMEZAWA Hiroyuki <kamezawa.hiroyu@jp.fujitsu.com>

---

---

Subject: Re: [PATCH v8 5/9] per-netns ipv4 sysctl\_tcp\_mem  
Posted by [KAMEZAWA Hiroyuki](#) on Fri, 09 Dec 2011 01:54:10 GMT  
[View Forum Message](#) <> [Reply to Message](#)

---

On Mon, 5 Dec 2011 19:34:59 -0200



Glauber Costa <glommer@parallels.com> wrote:

> This patch allows each namespace to independently set up  
> its levels for tcp memory pressure thresholds. This patch  
> alone does not buy much: we need to make this values  
> per group of process somehow. This is achieved in the  
> patches that follows in this patchset.  
>  
> Signed-off-by: Glauber Costa <glommer@parallels.com>  
> CC: KAMEZAWA Hiroyuki <kamezawa.hiroyu@jp.fujitsu.com>  
> CC: David S. Miller <davem@davemloft.net>  
> CC: Eric W. Biederman <ebiederm@xmission.com>

Reviewed-by: KAMEZAWA Hiroyuki <kamezawa.hiroyu@jp.fujitsu.com>

---

---

Subject: Re: [PATCH v8 6/9] tcp buffer limitation: per-cgroup limit  
Posted by [KAMEZAWA Hiroyuki](#) on Fri, 09 Dec 2011 01:55:19 GMT  
[View Forum Message](#) <> [Reply to Message](#)

---

On Mon, 5 Dec 2011 19:35:00 -0200

Glauber Costa <glommer@parallels.com> wrote:

> This patch uses the "tcp.limit\_in\_bytes" field of the kmem\_cgroup to  
> effectively control the amount of kernel memory pinned by a cgroup.  
>  
> This value is ignored in the root cgroup, and in all others,  
> caps the value specified by the admin in the net namespaces'  
> view of tcp\_sysctl\_mem.  
>  
> If namespaces are being used, the admin is allowed to set a  
> value bigger than cgroup's maximum, the same way it is allowed  
> to set pretty much unlimited values in a real box.  
>  
> Signed-off-by: Glauber Costa <glommer@parallels.com>  
> CC: David S. Miller <davem@davemloft.net>  
> CC: Hiroyuki Kamezawa <kamezawa.hiroyu@jp.fujitsu.com>  
> CC: Eric W. Biederman <ebiederm@xmission.com>

Reviewed-by: KAMEZAWA Hiroyuki <kamezawa.hiroyu@jp.fujitsu.com>

---

---

Subject: Re: [PATCH v8 7/9] Display current tcp memory allocation in kmem cgroup  
Posted by [KAMEZAWA Hiroyuki](#) on Fri, 09 Dec 2011 01:56:47 GMT  
[View Forum Message](#) <> [Reply to Message](#)

---

On Mon, 5 Dec 2011 19:35:01 -0200

Glauber Costa <glommer@parallels.com> wrote:

> This patch introduces kmem.tcp.usage\_in\_bytes file, living in the  
> kmem\_cgroup filesystem. It is a simple read-only file that displays the  
> amount of kernel memory currently consumed by the cgroup.  
>  
> Signed-off-by: Glauber Costa <glommer@parallels.com>  
> Reviewed-by: Hiroyouki Kamezawa <kamezawa.hiroyu@jp.fujitsu.com>  
> CC: David S. Miller <davem@davemloft.net>  
> CC: Eric W. Biederman <ebiederm@xmission.com>

Reviewed-by: KAMEZAWA Hiroyuki <kamezawa.hiroyu@jp.fujitsu.com>

---

---

Subject: Re: [PATCH v8 9/9] Display maximum tcp memory allocation in kmem cgroup

Posted by [KAMEZAWA Hiroyuki](#) on Fri, 09 Dec 2011 01:57:42 GMT

[View Forum Message](#) <> [Reply to Message](#)

---

On Mon, 5 Dec 2011 19:35:03 -0200

Glauber Costa <glommer@parallels.com> wrote:

> This patch introduces kmem.tcp.max\_usage\_in\_bytes file, living in the  
> kmem\_cgroup filesystem. The root cgroup will display a value equal  
> to RESOURCE\_MAX. This is to avoid introducing any locking schemes in  
> the network paths when cgroups are not being actively used.  
>  
> All others, will see the maximum memory ever used by this cgroup.  
>  
> Signed-off-by: Glauber Costa <glommer@parallels.com>  
> Reviewed-by: Hiroyouki Kamezawa <kamezawa.hiroyu@jp.fujitsu.com>  
> CC: David S. Miller <davem@davemloft.net>  
> CC: Eric W. Biederman <ebiederm@xmission.com>

Reviewed-by: KAMEZAWA Hiroyuki <kamezawa.hiroyu@jp.fujitsu.com>

---

---

Subject: Re: [PATCH v8 3/9] socket: initial cgroup code.

Posted by [KAMEZAWA Hiroyuki](#) on Fri, 09 Dec 2011 02:05:50 GMT

[View Forum Message](#) <> [Reply to Message](#)

---

On Mon, 5 Dec 2011 19:34:57 -0200

Glauber Costa <glommer@parallels.com> wrote:

> The goal of this work is to move the memory pressure tcp  
> controls to a cgroup, instead of just relying on global  
> conditions.

>  
> To avoid excessive overhead in the network fast paths,  
> the code that accounts allocated memory to a cgroup is  
> hidden inside a static\_branch(). This branch is patched out  
> until the first non-root cgroup is created. So when nobody  
> is using cgroups, even if it is mounted, no significant performance  
> penalty should be seen.  
>  
> This patch handles the generic part of the code, and has nothing  
> tcp-specific.  
>  
> Signed-off-by: Glauber Costa <glommer@parallels.com>  
> CC: Kirill A. Shutemov <kirill@shutemov.name>  
> CC: KAMEZAWA Hiroyuki <kamezawa.hiroyu@jp.fujitsu.com>  
> CC: David S. Miller <davem@davemloft.net>  
> CC: Eric W. Biederman <ebiederm@xmission.com>  
> CC: Eric Dumazet <eric.dumazet@gmail.com>

I already replied Reviewed-by: but...

```
> +/* Writing them here to avoid exposing memcg's inner layout */
> +#ifdef CONFIG_CGROUP_MEM_RES_CTLR_KMEM
> +#ifdef CONFIG_INET
> +#include <net/sock.h>
> +
> +static bool mem_cgroup_is_root(struct mem_cgroup *memcg);
> +void sock_update_memcg(struct sock *sk)
> +{
> + /* A socket spends its whole life in the same cgroup */
> + if (sk->sk_cgrp) {
> + WARN_ON(1);
> + return;
> + }
> + if (static_branch(&memcg_socket_limit_enabled)) {
> + struct mem_cgroup *memcg;
> +
> + BUG_ON(!sk->sk_prot->proto_cgroup);
> +
> + rcu_read_lock();
> + memcg = mem_cgroup_from_task(current);
> + if (!mem_cgroup_is_root(memcg)) {
> + mem_cgroup_get(memcg);
> + sk->sk_cgrp = sk->sk_prot->proto_cgroup(memcg);
> + }
> + rcu_read_unlock();
> + }
> + }
```

Here, you do mem\_cgroup\_get() if !mem\_cgroup\_is\_root().

```
> +EXPORT_SYMBOL(sock_update_memcg);
> +
> +void sock_release_memcg(struct sock *sk)
> +{
> + if (static_branch(&memcg_socket_limit_enabled) && sk->sk_cgrp) {
> + struct mem_cgroup *memcg;
> + WARN_ON(!sk->sk_cgrp->memcg);
> + memcg = sk->sk_cgrp->memcg;
> + mem_cgroup_put(memcg);
> + }
> +}
>
```

You don't check !mem\_cgroup\_is\_root(). Hm, root memcg will not be freed by this ?

Thanks,  
-Kame

---

Subject: Re: [PATCH v8 1/9] Basic kernel memory functionality for the Memory Controller

Posted by [Glauber Costa](#) on Fri, 09 Dec 2011 12:40:00 GMT

[View Forum Message](#) <> [Reply to Message](#)

---

On 12/08/2011 11:21 PM, KAMEZAWA Hiroyuki wrote:

> On Mon, 5 Dec 2011 19:34:55 -0200

> Glauber Costa<glommer@parallels.com> wrote:

>

>> This patch lays down the foundation for the kernel memory component  
>> of the Memory Controller.

>>

>> As of today, I am only laying down the following files:

>>

>> \* memory.independent\_kmem\_limit

>> \* memory.kmem.limit\_in\_bytes (currently ignored)

>> \* memory.kmem.usage\_in\_bytes (always zero)

>>

>> Signed-off-by: Glauber Costa<glommer@parallels.com>

>> Reviewed-by: Kirill A. Shutemov<kirill@shutemov.name>

>> CC: Paul Menage<paul@paulmenage.org>

>> CC: Greg Thelen<gthelen@google.com>

>

> As I wrote, please CC Johannes and Michal Hocko for memcg related parts.

I forgot to add them to the patch itself, but they are in the CC list of the messages.

So they did get the mail.

> A few questions.

> ==

>> + val = !!val;

>> +

>> + if (parent&& parent->use\_hierarchy&&

>> + (val != parent->kmem\_independent\_accounting))

>> + return -EINVAL;

> ==

> Hm, why you check val != parent->kmem\_independent\_accounting ?

>

> if (parent&& parent->use\_hierarchy)

> return -EINVAL;

> ?

Because I thought that making sure that everybody in the chain is consistent, it will make things simpler for us. But I am happy to change that if you prefer.

> BTW, you didn't check this cgroup has children or not.

> I think

>

> if (this\_cgroup->use\_hierarchy&&

> !list\_empty(this\_cgroup->children))

> return -EINVAL;

>

Noted.

> ==

>> + /\*

>> + \* TODO: We need to handle the case in which we are doing

>> + \* independent kmem accounting as authorized by our parent,

>> + \* but then our parent changes its parameter.

>> + \*/

>> + cgroup\_lock();

>> + memcg->kmem\_independent\_accounting = val;

>> + cgroup\_unlock();

>

> Do we need cgroup\_lock() here ?

Well, I removed almost all instances of it from previous patches, so I guess this one can go as well.

Subject: Re: [PATCH v8 2/9] foundations of per-cgroup memory pressure controlling.

Posted by [Glauber Costa](#) on Fri, 09 Dec 2011 12:41:02 GMT

[View Forum Message](#) <> [Reply to Message](#)

---

On 12/08/2011 11:24 PM, KAMEZAWA Hiroyuki wrote:

> On Mon, 5 Dec 2011 19:34:56 -0200

> Glauber Costa<glommer@parallels.com> wrote:

>

>> This patch replaces all uses of struct sock fields' memory\_pressure,  
>> memory\_allocated, sockets\_allocated, and sysctl\_mem to accessor  
>> macros. Those macros can either receive a socket argument, or a mem\_cgroup  
>> argument, depending on the context they live in.

>>

>> Since we're only doing a macro wrapping here, no performance impact at all is  
>> expected in the case where we don't have cgroups disabled.

>>

>> Signed-off-by: Glauber Costa<glommer@parallels.com>

>> CC: David S. Miller<davem@davemloft.net>

>> CC: Hiroyuki Kamezawa<kamezawa.hiroyu@jp.fujitsu.com>

>> CC: Eric W. Biederman<ebiederm@xmission.com>

>> CC: Eric Dumazet<eric.dumazet@gmail.com>

>

> please get ack from network guys.

> from me.

> Reviewed-by: KAMEZAWA Hiroyuki<kamezawa.hiroyu@jp.fujitsu.com>

>

Ok.

I think that now with all the Reviewed-by:'s I will resend it as a  
Request for Inclusion for Dave (plus fixing the problem you noted in  
patch 1)

---

---

Subject: Re: [PATCH v8 3/9] socket: initial cgroup code.

Posted by [Glauber Costa](#) on Fri, 09 Dec 2011 12:43:00 GMT

[View Forum Message](#) <> [Reply to Message](#)

---

On 12/09/2011 12:05 AM, KAMEZAWA Hiroyuki wrote:

> On Mon, 5 Dec 2011 19:34:57 -0200

> Glauber Costa<glommer@parallels.com> wrote:

>

>> The goal of this work is to move the memory pressure tcp  
>> controls to a cgroup, instead of just relying on global  
>> conditions.

>>

>> To avoid excessive overhead in the network fast paths,  
>> the code that accounts allocated memory to a cgroup is

```
>> hidden inside a static_branch(). This branch is patched out
>> until the first non-root cgroup is created. So when nobody
>> is using cgroups, even if it is mounted, no significant performance
>> penalty should be seen.
>>
>> This patch handles the generic part of the code, and has nothing
>> tcp-specific.
>>
>> Signed-off-by: Glauber Costa<glommer@parallels.com>
>> CC: Kirill A. Shutemov<kirill@shutemov.name>
>> CC: KAMEZAWA Hiroyuki<kamezawa.hiroyu@jp.fujitsu.com>
>> CC: David S. Miller<davem@davemloft.net>
>> CC: Eric W. Biederman<ebiederm@xmission.com>
>> CC: Eric Dumazet<eric.dumazet@gmail.com>
>
> I already replied Reviewed-by: but...
Feel free. Reviews, the more, the merrier.
```

```
>
>
>> +/* Writing them here to avoid exposing memcg's inner layout */
>> +#ifdef CONFIG_CGROUP_MEM_RES_CTLR_KMEM
>> +#ifdef CONFIG_INET
>> +#include<net/socket.h>
>> +
>> +static bool mem_cgroup_is_root(struct mem_cgroup *memcg);
>> +void sock_update_memcg(struct sock *sk)
>> +{
>> + /* A socket spends its whole life in the same cgroup */
>> + if (sk->sk_cgrp) {
>> + WARN_ON(1);
>> + return;
>> + }
>> + if (static_branch(&memcg_socket_limit_enabled)) {
>> + struct mem_cgroup *memcg;
>> +
>> + BUG_ON(!sk->sk_prot->proto_cgroup);
>> +
>> + rcu_read_lock();
>> + memcg = mem_cgroup_from_task(current);
>> + if (!mem_cgroup_is_root(memcg)) {
>> + mem_cgroup_get(memcg);
>> + sk->sk_cgrp = sk->sk_prot->proto_cgroup(memcg);
>> + }
>> + rcu_read_unlock();
>> + }
>> +}
>> +}
>
>
```

```

> Here, you do mem_cgroup_get() if !mem_cgroup_is_root().
>
>
>> +EXPORT_SYMBOL(sock_update_memcg);
>> +
>> +void sock_release_memcg(struct sock *sk)
>> +{
>> + if (static_branch(&memcg_socket_limit_enabled)&& sk->sk_cgrp) {
>> + struct mem_cgroup *memcg;
>> + WARN_ON(!sk->sk_cgrp->memcg);
>> + memcg = sk->sk_cgrp->memcg;
>> + mem_cgroup_put(memcg);
>> + }
>> +}
>>
>
> You don't check !mem_cgroup_is_root(). Hm, root memcg will not be freed
> by this ?
>
No, I don't. But I check if sk->sk_cgrp is filled. So it is implied,
because we only fill in this value if !mem_cgroup_is_root().

```

---

Subject: Re: [PATCH v8 1/9] Basic kernel memory functionality for the Memory Controller

Posted by [Glauber Costa](#) on Fri, 09 Dec 2011 14:37:23 GMT

[View Forum Message](#) <> [Reply to Message](#)

---

On 12/08/2011 11:21 PM, KAMEZAWA Hiroyuki wrote:

```

> Hm, why you check val != parent->kmem_independent_accounting ?
>
> if (parent&& parent->use_hierarchy)
> return -EINVAL;
> ?
>
> BTW, you didn't check this cgroup has children or not.
> I think
>
> if (this_cgroup->use_hierarchy&&
>      !list_empty(this_cgroup->children))
> return -EINVAL;

```

How about this?

```
val = !!val;
```

```
/*
```

```
 * This follows the same hierarchy restrictions than
```



```

* mem_cgroup_hierarchy_write()
*/
if (!parent || !parent->use_hierarchy) {
    if (list_empty(&cgroup->children))
        memcg->kmem_independent_accounting = val;
    else
        return -EBUSY;
}
else
    return -EINVAL;

return 0;

```

---

Subject: RE: [PATCH v8 1/9] Basic kernel memory functionality for the Memory Controller

Posted by [David Laight](#) on Fri, 09 Dec 2011 14:44:45 GMT

[View Forum Message](#) <> [Reply to Message](#)

---

> How about this?

```

>
>     val = !!val;
>
>     /*
>      * This follows the same hierarchy restrictions than
>      * mem_cgroup_hierarchy_write()
>      */
>     if (!parent || !parent->use_hierarchy) {
>         if (list_empty(&cgroup->children))
>             memcg->kmem_independent_accounting = val;
>         else
>             return -EBUSY;
>     }
>     else
>         return -EINVAL;
>
>     return 0;

```

Inverting the tests gives easier to read code:

```

if (parent && parent->user_hierarchy)
    return -EINVAL;
if (!list_empty(&cgroup->children))
    return -EBUSY;
memcg->kmem_independent_accounting = val != 0;
return 0;

```

NFI about the logic...

On the face of it the tests don't seem related to each other  
or to the assignment!

David

---

---

Subject: Re: [PATCH v8 1/9] Basic kernel memory functionality for the Memory Controller

Posted by [Glauber Costa](#) on Fri, 09 Dec 2011 14:48:16 GMT

[View Forum Message](#) <> [Reply to Message](#)

---

On 12/09/2011 12:44 PM, David Laight wrote:

```
>
>> How about this?
>>
>>     val = !!val;
>>
>>     /*
>>      * This follows the same hierarchy restrictions than
>>      * mem_cgroup_hierarchy_write()
>>      */
>>     if (!parent || !parent->use_hierarchy) {
>>         if (list_empty(&cgroup->children))
>>             memcg->kmem_independent_accounting = val;
>>         else
>>             return -EBUSY;
>>     }
>>     else
>>         return -EINVAL;
>>
>>     return 0;
>
```

> Inverting the tests gives easier to read code:

```
>
> if (parent&& parent->user_hierarchy)
>     return -EINVAL;
> if (!list_empty(&cgroup->children))
>     return -EBUSY;
> memcg->kmem_independent_accounting = val != 0;
> return 0;
```

On the other hand, inconsistent with mem\_cgroup\_hierarchy\_write(), which applies the logic in the same way I did here.

> NFI about the logic...

> On the face of it the tests don't seem related to each other

> or to the assignment!

How so?

If parent's use\_hierarchy is set, we can't set this value (we need to have a parent for that to even matter).

We also can't set it if we already have any children - otherwise all the on-the-fly adjustments become hell-on-earth.

As for = val != 0, sorry, but I completely disagree this is easier than !!val. Not to mention the !!val notation is already pretty widespread in the kernel.

> David

>

>

>

>

> --

> To unsubscribe, send a message with 'unsubscribe linux-mm' in

> the body to majordomo@kvack.org. For more info on Linux MM,

> see: <http://www.linux-mm.org/> .

> Fight unfair telecom internet charges in Canada: sign <http://stopthemeteter.ca/>

> Don't email:<a href=mailto:dont@kvack.org"> email@kvack.org</a>

---

---

Subject: Re: [PATCH v8 3/9] socket: initial cgroup code.

Posted by [KAMEZAWA Hiroyuki](#) on Mon, 12 Dec 2011 00:33:13 GMT

[View Forum Message](#) <> [Reply to Message](#)

---

On Fri, 9 Dec 2011 10:43:00 -0200

Glauber Costa <[glommer@parallels.com](mailto:glommer@parallels.com)> wrote:

> On 12/09/2011 12:05 AM, KAMEZAWA Hiroyuki wrote:

> > On Mon, 5 Dec 2011 19:34:57 -0200

> > Glauber Costa<[glommer@parallels.com](mailto:glommer@parallels.com)> wrote:

> >

> >> The goal of this work is to move the memory pressure tcp

> >> controls to a cgroup, instead of just relying on global

> >> conditions.

> >>

> >> To avoid excessive overhead in the network fast paths,

> >> the code that accounts allocated memory to a cgroup is

> >> hidden inside a static\_branch(). This branch is patched out

> >> until the first non-root cgroup is created. So when nobody

> >> is using cgroups, even if it is mounted, no significant performance

> >> penalty should be seen.

> >>

> >> This patch handles the generic part of the code, and has nothing

```

> >> tcp-specific.
> >>
> >> Signed-off-by: Glauber Costa<glommer@parallels.com>
> >> CC: Kirill A. Shutemov<kirill@shutemov.name>
> >> CC: KAMEZAWA Hiroyuki<kamezawa.hiroyu@jp.fujitsu.com>
> >> CC: David S. Miller<davem@davemloft.net>
> >> CC: Eric W. Biederman<ebiederm@xmission.com>
> >> CC: Eric Dumazet<eric.dumazet@gmail.com>
> >
> > I already replied Reviewed-by: but...
> Feel free. Reviews, the more, the merrier.
>
> >
> >
> >> +/* Writing them here to avoid exposing memcg's inner layout */
> >> +#ifdef CONFIG_CGROUP_MEM_RES_CTLR_KMEM
> >> +#ifdef CONFIG_INET
> >> +#include<net/sock.h>
> >> +
> >> +static bool mem_cgroup_is_root(struct mem_cgroup *memcg);
> >> +void sock_update_memcg(struct sock *sk)
> >> +{
> >> + /* A socket spends its whole life in the same cgroup */
> >> + if (sk->sk_cgrp) {
> >> + WARN_ON(1);
> >> + return;
> >> + }
> >> + if (static_branch(&memcg_socket_limit_enabled)) {
> >> + struct mem_cgroup *memcg;
> >> +
> >> + BUG_ON(!sk->sk_prot->proto_cgroup);
> >> +
> >> + rcu_read_lock();
> >> + memcg = mem_cgroup_from_task(current);
> >> + if (!mem_cgroup_is_root(memcg)) {
> >> + mem_cgroup_get(memcg);
> >> + sk->sk_cgrp = sk->sk_prot->proto_cgroup(memcg);
> >> + }
> >> + rcu_read_unlock();
> >> + }
> >> +}
> >
> > Here, you do mem_cgroup_get() if !mem_cgroup_is_root().
> >
> >
> >> +EXPORT_SYMBOL(sock_update_memcg);
> >> +
> >> +void sock_release_memcg(struct sock *sk)

```

```

> >> +{
> >> + if (static_branch(&memcg_socket_limit_enabled)&& sk->sk_cgrp) {
> >> + struct mem_cgroup *memcg;
> >> + WARN_ON(!sk->sk_cgrp->memcg);
> >> + memcg = sk->sk_cgrp->memcg;
> >> + mem_cgroup_put(memcg);
> >> + }
> >> +}
> >>
> >
> > You don't check !mem_cgroup_is_root(). Hm, root memcg will not be freed
> > by this ?
> >
> > No, I don't. But I check if sk->sk_cgrp is filled. So it is implied,
> > because we only fill in this value if !mem_cgroup_is_root().

```

Ah, ok. thank you.  
-Kame

---

Subject: Re: [PATCH v8 1/9] Basic kernel memory functionality for the Memory Controller

Posted by [KAMEZAWA Hiroyuki](#) on Mon, 12 Dec 2011 00:34:48 GMT

[View Forum Message](#) <> [Reply to Message](#)

---

On Fri, 9 Dec 2011 12:37:23 -0200

Glauber Costa <glommer@parallels.com> wrote:

```

> On 12/08/2011 11:21 PM, KAMEZAWA Hiroyuki wrote:
> > Hm, why you check val != parent->kmem_independent_accounting ?
> >
> > if (parent&& parent->use_hierarchy)
> > return -EINVAL;
> > ?
> >
> > BTW, you didn't check this cgroup has children or not.
> > I think
> >
> > if (this_cgroup->use_hierarchy&&
> >      !list_empty(this_cgroup->children))
> > return -EINVAL;
> >
> > How about this?
> >
> >     val = !!val;
> >
> >     /*
> >      * This follows the same hierarchy restrictions than

```

```
> * mem_cgroup_hierarchy_write()
> */
> if (!parent || !parent->use_hierarchy) {
>     if (list_empty(&cgroup->children))
>         memcg->kmem_independent_accounting = val;
>     else
>         return -EBUSY;
> }
> else
>     return -EINVAL;
>
> return 0;
>
```

seems good to me.

Thanks,  
-Kame

---