Semantics and Behavior of Local Atomic Operations

Author: Mathieu Desnoyers

This document explains the purpose of the local atomic operations, how to implement them for any given architecture and shows how they can be used properly. It also stresses on the precautions that must be taken when reading those local variables across CPUs when the order of memory writes matters.

Note

Note that <code>local_t</code> based operations are not recommended for general kernel use. Please use the <code>this_cpu</code> operations instead unless there is really a special purpose. Most uses of <code>local_t</code> in the kernel have been replaced by <code>this_cpu</code> operations. <code>this_cpu</code> operations combine the relocation with the <code>local_t</code> like semantics in a single instruction and yield more compact and faster executing code.

Purpose of local atomic operations

Local atomic operations are meant to provide fast and highly reentrant per CPU counters. They minimize the performance cost of standard atomic operations by removing the LOCK prefix and memory barriers normally required to synchronize across CPUs.

Having fast per CPU atomic counters is interesting in many cases: it does not require disabling interrupts to protect from interrupt handlers and it permits coherent counters in NMI handlers. It is especially useful for tracing purposes and for various performance monitoring counters.

Local atomic operations only guarantee variable modification atomicity wrt the CPU which owns the data. Therefore, care must taken to make sure that only one CPU writes to the $local_t$ data. This is done by using per cpu data and making sure that we modify it from within a preemption safe context. It is however permitted to read $local_t$ data from any CPU: it will then appear to be written out of order wrt other memory writes by the owner CPU.

Implementation for a given architecture

It can be done by slightly modifying the standard atomic operations: only their UP variant must be kept. It typically means removing LOCK prefix (on i386 and x86_64) and any SMP synchronization barrier. If the architecture does not have a different behavior between SMP and UP, including asm-generic/local.h in your architecture's local.h is sufficient.

The $local_t$ type is defined as an opaque signed long by embedding an $atomic_long_t$ inside a structure. This is made so a cast from this type to a long fails. The definition looks like:

```
typedef struct { atomic_long_t a; } local_t;
```

Rules to follow when using local atomic operations

- Variables touched by local ops must be per cpu variables.
- Only the CPU owner of these variables must write to them.
- This CPU can use local ops from any context (process, irq, softirq, nmi, ...) to update its local t variables.
- Preemption (or interrupts) must be disabled when using local ops in process context to make sure the process won't be migrated to a different CPU between getting the per-cpu variable and doing the actual local op.
- When using local ops in interrupt context, no special care must be taken on a mainline kernel, since they will run on the local CPU with preemption already disabled. I suggest, however, to explicitly disable preemption anyway to make sure it will still work correctly on -rt kernels.
- Reading the local cpu variable will provide the current copy of the variable.
- Reads of these variables can be done from any CPU, because updates to "long", aligned, variables are always atomic. Since
 no memory synchronization is done by the writer CPU, an outdated copy of the variable can be read when reading some
 other cpu's variables.

How to use local atomic operations

```
#include <linux/percpu.h>
#include <asm/local.h>
static DEFINE PER CPU(local t, counters) = LOCAL INIT(0);
```

Counting

Counting is done on all the bits of a signed long.

In preemptible context, use get cpu var () and put cpu var () around local atomic operations: it makes sure that preemption is

disabled around write access to the per cpu variable. For instance:

```
local_inc(&get_cpu_var(counters));
put cpu var(counters);
```

If you are already in a preemption-safe context, you can use this cpu ptr() instead:

```
local_inc(this_cpu_ptr(&counters));
```

Reading the counters

Those local counters can be read from foreign CPUs to sum the count. Note that the data seen by local_read across CPUs must be considered to be out of order relatively to other memory writes happening on the CPU that owns the data:

If you want to use a remote local_read to synchronize access to a resource between CPUs, explicit $smp_wmb()$ and $smp_rmb()$ memory barriers must be used respectively on the writer and the reader CPUs. It would be the case if you use the local_t variable as a counter of bytes written in a buffer: there should be a $smp_wmb()$ between the buffer write and the counter increment and also a $smp_rmb()$ between the counter read and the buffer read.

Here is a sample module which implements a basic per cpu counter using local.h:

```
/* test-local.c
 * Sample module for local.h usage.
#include <asm/local.h>
#include <linux/module.h>
#include <linux/timer.h>
static DEFINE PER CPU(local t, counters) = LOCAL INIT(0);
static struct timer list test timer;
/* IPI called on each CPU. */
static void test_each(void *info)
        /* Increment the counter from a non preemptible context */
        printk("Increment on cpu %d\n", smp_processor_id());
        local_inc(this_cpu_ptr(&counters));
        /* This is what incrementing the variable would look like within a
        * preemptible context (it disables preemption) :
         * local inc(&get_cpu_var(counters));
         * put_cpu_var(counters);
}
static void do test timer(unsigned long data)
        int cpu;
        /* Increment the counters */
        on each cpu(test each, NULL, 1);
        /* Read all the counters */
        printk("Counters read from CPU %d\n", smp processor id());
        for each online cpu(cpu) {
                printk("Read : CPU %d, count %ld\n", cpu,
                        local read(&per cpu(counters, cpu)));
       mod_timer(&test_timer, jiffies + 1000);
static int init test init(void)
        /* initialize the timer that will increment the counter */
        timer_setup(&test_timer, do_test_timer, 0);
        mod_timer(&test_timer, jiffies + 1);
        return 0;
static void __exit test_exit(void)
```

```
del_timer_sync(&test_timer);
}
module_init(test_init);
module_exit(test_exit);

MODULE_LICENSE("GPL");
MODULE_AUTHOR("Mathieu Desnoyers");
MODULE_DESCRIPTION("Local Atomic Ops");
```