



The Labs and Linux

Revision 1.3

Larry McVoy

BitMover, Inc.
San Francisco, California
415-821-5758
lm@bitmover.com

ABSTRACT

The computing world is evolving at a rapid pace. While most segments of the computer market are growing, the supercomputer area has not kept up with the expansion of other areas in the market. In fact, the supercomputer segment is in danger of being relegated to a small niche player.

The Labs have a vested interest in the health of the supercomputer market, being one of the primary users of supercomputer technology. The Labs have very specific needs which are not currently being well addressed by either the commercial vendors or the free software world.

In this document, we will contrast the Labs needs versus the needs of the more general computing market. The role of Linux with respect to supercomputing will be examined, with special attention being paid to Linux' clustering plans. We will consider several options open to the Labs. Finally, we will suggest a moderate approach to solving the computing problems facing the Labs. We believe that it is possible for the Labs to get what they want with a fairly small investment and some innovative approaches to the problem.

The approach we suggest involves Linux clusters, so the second half of this document describes what we mean by a cluster, why this view is important, and then finishes by explaining why this sort of cluster is important to the Labs. It may seem obvious that the Labs want clusters, but it is not immediately obvious why the Labs would want the sort of clusters described here.

6.1. Why SMP clusters?

A good question is: why should this sort of technology ever exist? We already have several operating systems which scale up to 64 processors and at least one example, IRIX, which scales up to 256 processors. Why bother with a different way to solve the same problem?

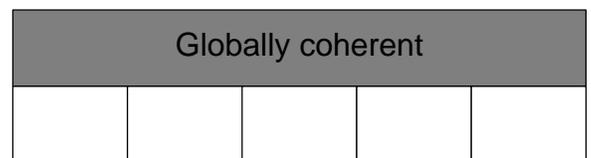
Consider the different demands placed on a 2 processor SMP OS and a 256 processor SMP OS. They both offer the same abstractions but the larger machine needs an operating system which can have 256 processors all doing something useful - in the kernel - at the same time. That means the 256 processor OS is fair more fine grained in its locking approach than the 2 processor case. The problems associated with locking are discussed in another paper, suffice it to say that asking one operating system to handle uniprocessor and 2-256 SMP operating system tasks is an unreasonable request.

There is another consideration, which is Linux specific. Linux today is not a fine grained threaded OS, it scales only to a few processors. That means that Linux is still a relatively simple OS by today's standards. That's a good thing. Scaling Linux to 100 or 200 processors would turn Linux into something that was no longer the lightweight, high performance OS we know on small machines. It would trade low end performance for high end performance, which is not a good trade off.

None the less, not scaling is not an option. There are problems which are larger than 4 processor systems. Linux has to be able to solve those problems if it is to be taken seriously in the enterprise market. The pressure to scale up is already here and is being felt on a daily basis by the Linux developers. There has to be a scaling answer.

SMP OS clusters are a less invasive approach to scaling and could be applied to Linux. While the virtual process changes are quite invasive, the rest of the changes are not - they tend to fit nicely under existing abstractions such as the file system layer. The problem is nicely partitionable as well - parallel projects could be started on the process model, the load balancing, the file system coherency, etc. This approach has the built in advantage that it is inherently more scalable; consider scaling up a run queue on a 1024 processor SMP system and then consider scaling it up on an SMP OS cluster - the partitioning needed for scaling is already done.

The final plug for this model is this: we have knowledge that at least two >\$10B/year hardware vendors have designed and are building hardware with features specifically designed to support this model. In particular, the hardware has support for multiple coherency domains, so each OS can see two kinds of memory: locally coherent and globally coherent. You can then have a view of memory like this:



4xCPU 4xCPU 4xCPU 4xCPU 4xCPU

where each lower box is memory which is coherent for a single 4 processor node running one OS image, and the upper box is a large portion of memory which is globally coherent.

The reason for building such an architecture is that SMP doesn't really scale - it gets harder and harder to build systems with globally coherent memory. Once you realize that you are going to essentially partition all of the important data structures in the OS in order to scale up, having those partitioned data structures in localized

