# A NEW SCALABLE MULTI-NODE
# EVENT LOGGING SYSTEM FOR BABAR

JAMES A. HAMILTON

*Stanford Linear Accelerator Center, 2575 Sand Hill Road*
*Menlo Park, California 94025, U.S.A.*

STEFFEN LUITZ

*Stanford Linear Accelerator Center, 2575 Sand Hill Road*
*Menlo Park, California 94025, U.S.A.*

For the BaBar Computing Group

The BaBar experiment is currently operating near the rate limit of its ability to log event data to disk and tape using the existing hardware and software systems. Consequently we have chosen to design and implement a new system for logging event data. The new system is designed to be scalable, so that the data rate can be increased by adding systems at one of three levels. It also has the property that data can be logged at almost unlimited burst rates without introducing dead time. The key to these features lies in the use of many nodes within the level three trigger system of BaBar. This allows the events to first be logged to local disks within the trigger system, and then later to be merged to any of multiple merge servers in non-real-time.

*Keywords*: Events; logging; scalability.

## 1. Introduction

The BaBar experiment has been in operation at SLAC for about five years, since 1999. It takes data continuously, with the exception of breaks for maintenance and upgrades. By the end of 2003, the rate of events logged to disk and tape had reached 300 per second. With a mean event size of about 30 kilobytes, this amounts to about 9 megabytes/second. The event logging system which was in place at the time was capable of only about 10 megabytes/second. By 2007, we expect the data rate to be 500 kilohertz at 75 kilobytes per event. So it was clearly necessary to replace or upgrade this system.

We chose to redesign the event logging system. Along with meeting the lifetime data rate expectations for the experiment, this redesign had three additional goals.
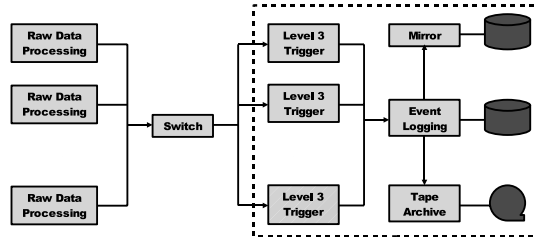
Fig. 1.   The existing BaBar data acquisition system, The portion to be redesigned is outlined.

First, we wanted the system to be fault tolerant, so that we could continue taking data if any single part of the system failed. Second, we wanted it to be scalable, so that it could grow incrementally to accommodate higher data rates. Finally, we wanted to occasionally take data at very high rates that would exceed the average capability of the logging system, and yet to do this without losing data or causing dead time at the detector.

## 2.  The Old Logging System

Figure 1 shows a simplified version of the data acquisition system as it existed prior to this work. Raw data is collected at the detector based on the level one trigger system. These events flow through a network switch, and each event is passed to one of the several level three trigger systems. There it is analyzed and filtered so that the event rate can be reduced by approximately a factor of ten.

These events were then sent over 100 megabit links to the logging server. For convenience, the data stream, even though it is essentially continuous, is divided into units called runs. Each run is approximately one hour's worth of data. The run is the unit of data stored in a file by the logging system. After the completion of a run, the file is then sent over a gigabit ethernet to tape for permanent storage. An additional copy of the data is made to a mirror system for use by subsystems that may want immediate access to the data.

To support these three data streams in the old system would have required a throughput of 113 megabytes/sec. This is beyond the capability of reasonably priced current systems.

## 3.  The New Logging Manager Architecture

Figure 2 shows a typical setup for logging events in the new architecture. The most obvious feature of this architecture, as mentioned above, is that the real-time event logging is to the local disks on the level three trigger systems. Later, the data is merged onto merge server disks in non-real-time.

If we assume a throughput of 50 megabytes/second to local disks, we can log about 20,000 hertz of 75 kilobyte events. This is well in excess of anything we expect the level one trigger to produce.
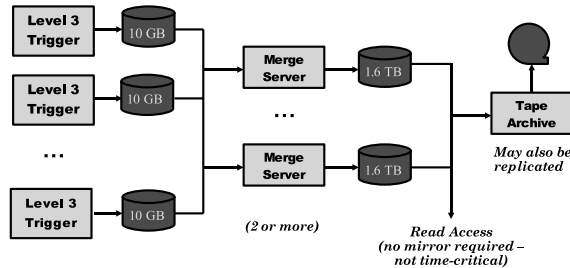
Fig. 2.   The components and data paths for logging events in the new architecture.

### 3.1. *The Merge Subsystem*

We have measured the merge rate using the new system at 80 megabytes/second. Even with only one server, this rate is eight times the previous capability. We have installed a second server for redundancy. If there is a problem with either of these server, one can be taken out of service indefinitely without compromising the ability to take data.

We must still transfer the data to the tape system, but now this transfer *follows* the merge transfer so there is no interference. The merge server remains only about 25 utilized, so we can dispense with the mirror server and allow direct nfs access to the logging servers.

## 4. Failure Analysis

This new system is fault tolerant in the sense that we can continue taking data after any single point of failure in the event logging system, because all data disks are either fully raided or mirrored. But it is also important to see that the logging system never loses data that has already been taken. We do this by recording a merge only when all data has been written to disk on a merge server. Incomplete or failed merges will be retried from the beginning. All system components keep sufficient state information to determine the current location of all event data for all recent runs.

## 5. Conclusions

We have developed a new event logging system for BaBar. In addition to providing the higher throughput required for the expected future data rates, it meets several important additional goals. First, it is scalable. Second, it supports high burst-rate data taking without introducing dead time. Finally, it is fault tolerant in terms of both allowing data-taking to continue while repairs are underway; and also in terms of preventing the loss of any data already collected.

The system has been in place since July, 2004, and has proved very reliable, and has shown that it meets these goals in practice as well as in theory.