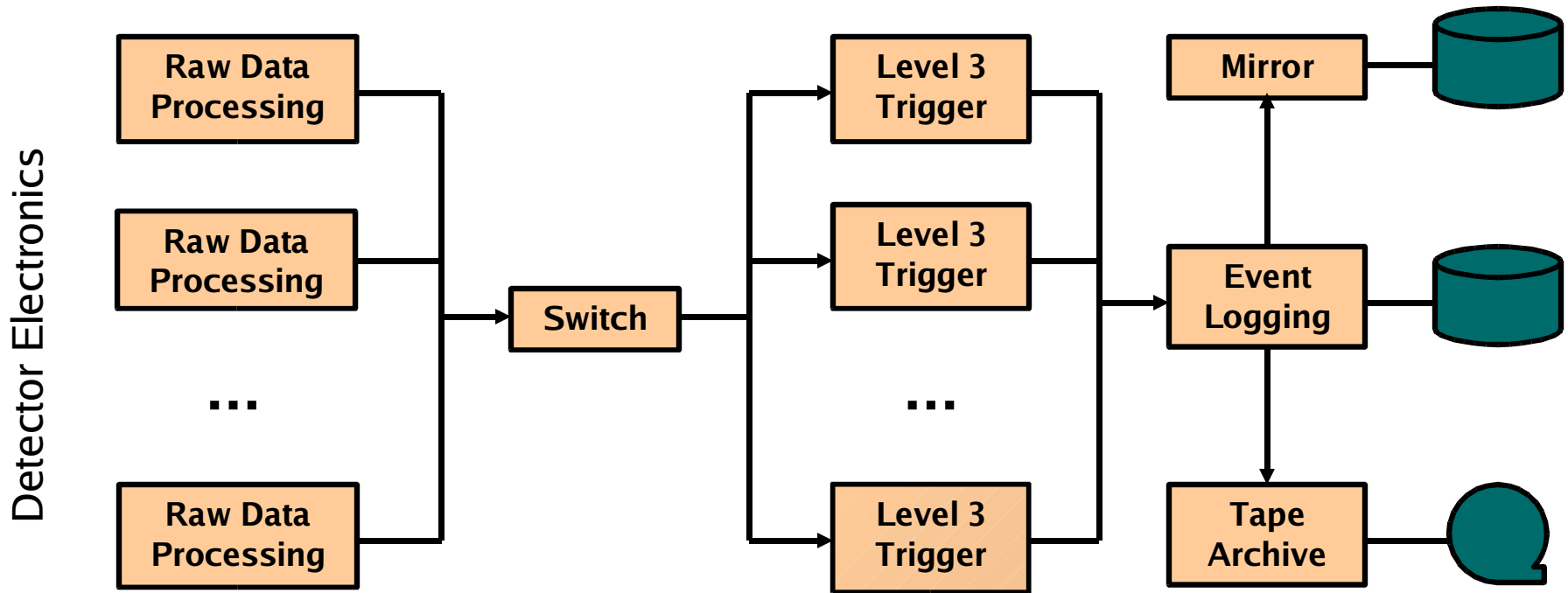


A New Scalable Multi-Node Event Logging System for BaBar

James A. Hamilton
Steffen Luitz

For the BaBar Computing Group

Original Structure





Performance Requirements

- **Current**
 - ✓ *30 kilobyte events X 300 events/sec. = 9 megabytes/sec.*
 - ✓ *Server limited to ~ 10 megabytes/sec.*
- **Projected by 2007**
 - ✓ *75 kilobyte events X 500 events/sec. = 37.5 megabytes/sec.*



Options

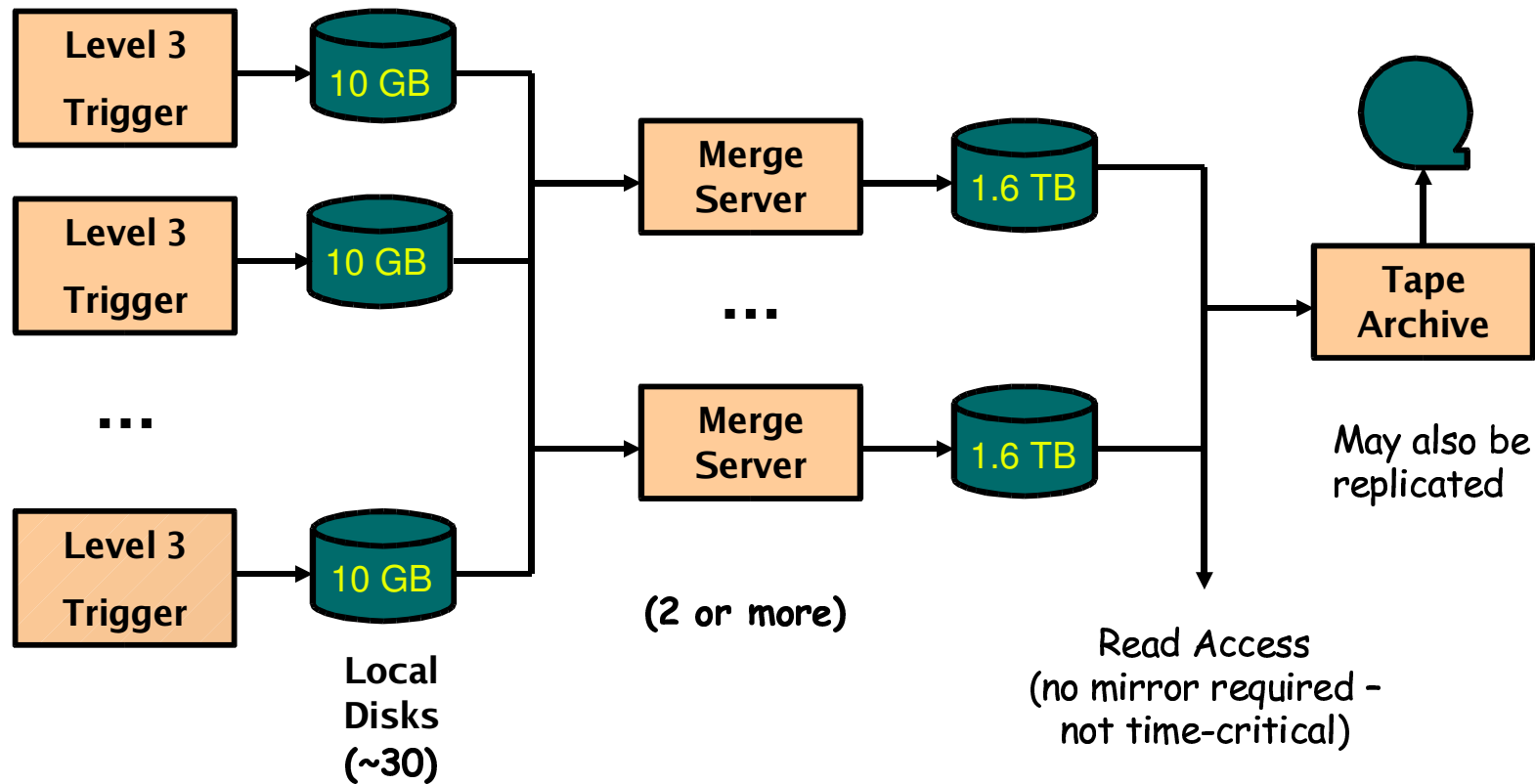
- Could do a simple HW/SW upgrade
- But a new architecture can also get:
 - *Scalability in case we guess wrong*
 - *Fault tolerance (RAID array not enough because of long disk reconstruction time)*
 - *High burst rate to allow passthrough runs*
- So we chose the new architecture



New Architecture

- Take advantage of multiple level 3 systems (currently 30)
- Write data in parallel to local disks on level 3 systems
- Merge to file server at end of run
 - *BaBar B-factory runs continuously with ~1 hour runs*
 - *Uses trickle injection, so data-taking never stops*
- Multiple file servers for finished runs

Data Paths





Server Performance

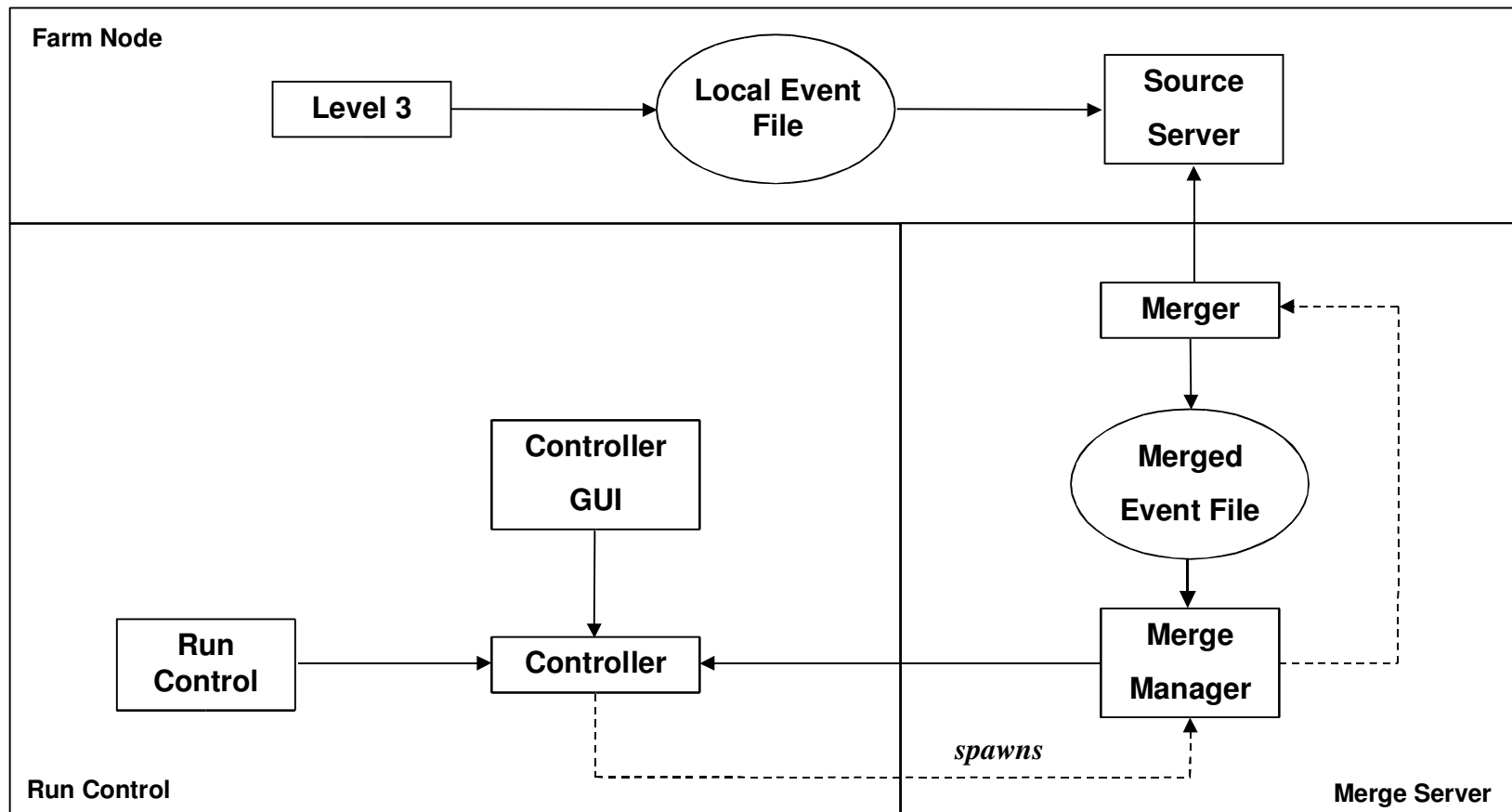
- Each merge server is equipped with:
 - *1.6 terabyte raid disk array at ~100 megabytes/sec.*
 - *Gigabit ethernet*
 - *Dual 3GHz P4 running Linux & XFS*
- Merge operations average 80 megabytes/sec. transfer rate
 - *Over 1000 events/sec. At 75 kilobytes per event*
- Use second server for fault tolerance



Level 3 System Performance

- Each level 3 system has:
 - *Two 18 gigabyte mirrored disks at ~50 megabytes/sec.*
 - Leaves ~10 gigabytes for event data – enough for ~10 hours at current rates
 - Need larger disks here already
 - *Gigabit ethernet to dataflow*
 - *100 megabit ethernet to merge servers*
- Each merger takes 2.6 megabytes/sec. (80/30)
 - *5 % of disk bandwidth*
 - *21 % of network bandwidth*
- ~20,000 events/sec. burst, at 75 kilobytes/event
 - *Have done 3,000 events/sec. in production*

Component Architecture



Controller GUI

Run#	State	Server	Date-Time
0050506	Waiting	NONE	07/27-17:39
0050505	Complete	bfo-log101	07/27-16:44
0050504	Complete	bfo-log100	07/27-15:49
0050503	Complete	bfo-log101	07/27-14:56
0050502	Complete	bfo-log100	07/27-14:00
0050501	Complete	bfo-log101	07/27-13:05
0050500	Complete	bfo-log100	07/27-12:10
0050499	Complete	bfo-log101	07/27-11:15
0050498	Complete	bfo-log100	07/27-10:31
0050497	Complete	bfo-log101	07/27-09:35
0050496	Complete	bfo-log100	07/27-08:40
0050495	Waiting	NONE	07/27-07:44
0050494	Complete	bfo-log100	07/27-06:49

Server	Run
bbr-nfs102	OFFLINE
bbt-srv100	FULL
bfo-log100	OFFLINE
bfo-log101	OFFLINE

ERROR No server available



Failure Modes

- Loss of network connection
- Unrecoverable single disk error
- Processor failure
- Component crash or restart
- Must recover from these without loss of data



Merge Failure Recovery

- Repair
 - *Replace failing drive (recover via RAID) or processor*
 - *Continue taking data during repair*
- Retry
 - *Source data never deleted until merge succeeds*
 - *Manual and automatic*
- Scavenge
 - *Controller queries source servers and merge manager*
- Redirect
 - *Sources can be renamed or omitted*



Conclusions

- We have developed a new event logging system that is:
 - *Scalable – meets projected throughput requirements for lifetime of BaBar*
 - *Fault tolerant*
 - *Supports high burst-rate data taking*
- In successful production use for last two weeks of Run 4. Will continue in Run 5.



Scalability Timelines

