

# Bandwidth needs in HPC taking into account link redundancy

Petar Pepeljugoski, Paul Coteus

IBM

11/2011

# HPC/server environment

- The performance improvement trend continues unabated
- Government target is an exascale machine by 2019
  - Can't be achieved in one step
  - By 2016 we will be seeing 400 PF machines
- Performance improvement is more than two orders of magnitude from today's fastest machines
  - When we argued for 100Gb/s Ethernet, performance target was 1 PF, need much MOORE bandwidth

# How to achieve higher link bandwidth

- Can try to increase line data rate, but inconceivable to achieve two orders of magnitude improvement
- Can change architecture, to lower bandwidth requirement – but only to a point
- Increase parallelism – apply to optical (SMF and MMF links) and electrical links (if they are still around) for 400 PF machine,
  - 400 Gb/s (16 send, 16 receive at 25 Gb/s) is likely
- Whatever approach HPC needs lot of links
- With current FIT rates there will be lot of fails, so need spares
- Need to think about packaging and replacement strategy
- Industry has to decide which is more economical
  - Use fail in place philosophy. If lane in a cable fails, work around it.
  - Make cable separable into active transceiver and passive cable, so cable can remain untouched while failed transceiver is replaced (current approach), add redundant channel
  - Fuse transceiver and cable (active optical cable) and add redundant channel

# Failure Analysis

- Current links have no lane redundancy
  - Increasing parallelism makes them more vulnerable to failure
- Typical fit rate is 10 FIT per lane (16 wide link x 2 is 320 FIT)
  - 30000 link (not unreasonable number in HPC/server environment) would mean 84 single fiber fails per year
- This is unacceptable, need link redundancy (one extra lane) to reduce fit rate to 10 FIT
- Spare fiber may have other uses, like in power reduction schemes

# Conclusion

- We need at least 400 Gb/s in future HPC/server environments
- Need to invest in extra bandwidth for lane redundancy