**Leigh, B. Kevin, "Design and Analysis of Network and IO Consolidations in a General-Purpose Infrastructure"**

*Advisor: Dr. J. Subhlok*

Our goal for this dissertation is to prove our hypothesis that IO consolidation is a cost efficient future server strategy. IO consolidation is to replace dedicated network interface devices (in servers) and corresponding network switches, with lower cost IO switches and fewer shared network interface devices. Our key contributions include, (a) a novel general-purpose infrastructure that seamlessly supports traditional and future interconnect protocols in server blade enclosures; and (b) a unique method, including new cost and performance models, to evaluate efficiency of IO consolidation with respect to traditional networks, leveraging our general-purpose infrastructure as the framework. The first part of the dissertation focuses on the design of a general-purpose infrastructure (GPI). We combine technology trends and our insights into several fabrics' physical layer similarities. We then architected an interconnect infrastructure by volume-metric positioning of switch bays and adaptive grouping of signal lanes. The resulting infrastructure is general-purpose enough to support different fabric protocols. The second part of the dissertation evaluates performance and cost efficiencies of IO consolidation solutions, with respect to a network consolidation solution, within a GPI. We choose two industry-standard protocols for the respective solutions – the PCI Express as our IO consolidation protocol since it has become the IO interface in high volume computers, and the 10Gbit Ethernet as our reference network consolidation protocol since it has become the de facto network used in replace of other networks. Because of the newness of these two protocols there are not enough full systems to easily compare the two corresponding consolidation methods. We solve this problem by deriving our own performance metrics and cost models based on our insights into the characteristics of the protocols, the constraints of the GPI, and the characteristics of datacenter applications. We use a hybrid method consisting of hardware emulators, software simulators and hardware prototypes. Our main results prove that IO consolidation can be performance and cost efficient within the GPI architecture, but require careful tradeoffs. Our results also show that smaller IO fabrics are compelling for their dramatic cost savings and scalability, while maintaining high enough performance with well-balanced bandwidth utilization throughout the infrastructure.