

**FY16 Alternatives Analysis
for the
Lattice QCD Computing Project Extension II
(LQCD-ext II)**

Operated at
Brookhaven National Laboratory
Fermi National Accelerator Laboratory
Thomas Jefferson National Accelerator Facility

for the
U.S. Department of Energy
Office of Science
Offices of High Energy and Nuclear Physics

Version 0.5

Revision Date
July 21, 2015

PREPARED BY:
Chip Watson, Jefferson Lab

CONCURRENCE:

William N. Boroski
LQCD-ext Contract Project Manager

Date

**Lattice QCD Computing Project Extension (LQCD-ext)
Change Log: Alternatives Analysis for FY16 Procurement**

Revision No.	Description	Effective Date
0.1	Document created from FY14 document.	July 12, 2015
0.2	Added additional Intel Omni-Path information released at ISC.	July 13, 2015
0.3	Updated with Intel Haswell data, improved data from DH on Pi0g, additional Pascal NVlink data	July 14, 2015
0.4	Increased uncertainty in cost and in possible higher performance per socket of KNL alternative; expanded discussion	July 20, 2015
0.5	Added more details on existing KNL software, made other minor adjustments and editorial fixes	July 21, 2015

Table of Contents

1	Introduction.....	1
2	FY16 Goals.....	1
3	Hardware Options.....	2
4	Alternatives.....	7
4.1	Alternative 1: A Xeon Phi / KNL cluster released to production by Sept 30, 2016.	8
4.2	Alternative 2: A 50% - 50% (by budget) mixture of conventional and GPU-accelerated clusters released to production by Sept 30, 2016	8
4.3	Alternative 3: A pure GPU-accelerated cluster released to production by Sept 30, 2016 ...	9
4.4	Alternative 4: A conventional cluster released to production by July 1, 2016	9
4.5	Alternative 5: Expand the half-rack of BG/Q deployed in Q1 2013 to a full rack, and deploy a small GPU cluster in FY16.....	10
4.6	Alternative 6: Status Quo (no additional deployment in FY16)	10
4.7	Other Alternatives	11
5	Discussion.....	11
6	Conclusion	12

1 Introduction

This document presents the analysis of FY16 alternatives for obtaining the computational capacity needed for the US Lattice QCD effort within High Energy Physics (HEP) and Nuclear Physics (NP) by the SC Lattice QCD Computing Extension Project (LQCD-ext). This analysis is updated at least annually to capture decisions taken during the life of the project, and to examine options for the next year. The technical managers of the project are also continuously tracking market developments through interactions with computer and chip vendors, through trade journals and online resources, and through computing conferences. This tracking allows unexpected changes to be incorporated into the project execution in a timely fashion.

Alternatives herein are constrained to approximately fit within the current budget guidance of the project, in particular ~\$0.86M for computing and file server procurements in FY 2016. This constraint provides adequate funding to meet the basic requirements of the field for enhanced computational capacity, under the assumption of expanding resources at ANL and ORNL already planned by the Office of Science (SC), and under the assumption that a reasonable fraction of those resources are ultimately allocated to Lattice QCD.

All alternatives assume the continued operation of the existing resources from the FY09-FY13 LQCD Facilities Projects until those resources reach end of life, i.e., until each resource is no longer cost effective to operate, typically about 5 years. At present these resources constitute an aggregate conventional x86 resource of about 70 teraflop/s sustained on LQCD benchmarks, 650 GPUs with an effective capacity of about 100 teraflop/s sustained, and a half rack of BG/Q with a performance of about 16 teraflops, for a total of around 186 teraflops sustained. The aggregate project cost of operating these existing systems in FY2016 is approximately \$1.72M (for the three sites combined). Replacing and running the flexible computational capacity represented by these existing resources cannot be done for less than its current operating cost.

In FY16, viable hardware options are a conventional Infiniband cluster, a GPU accelerated cluster, a Xeon Phi (Knights Landing) cluster, or some combination of these. Conventional clusters can run codes for all actions of interest to USQCD. Optimized multi-GPU codes are available for the HISQ, Wilson, clover, and twisted mass actions, with code for the DWF action and for multi-grid under development. Optimized inverter software on Xeon Phi / Knights Corner is available for Wilson, clover (USQCD SciDAC software) HISQ (international collaboration) and FGMRES-DR (international, optimization helped by JLab's QPhiX code generator). Several additional actions used in USQCD software, including multi-grid and DWF, are currently also under development for Xeon Phi / Knights Landing.

2 FY16 Goals

The project baseline calls for deployment in FY16 of 49 TF sustained performance, based upon extrapolations of price performance of Intel x86 cores and NVIDIA Tesla GPUs, and an assumption of using 60% of the compute budget for conventional x86 nodes, and 40% for GPU accelerated nodes.

The choice of 60:40 split was driven by the recognition that there is still substantial software that cannot (yet) exploit GPUs, so that 100% of the funds for GPUs would not be the best choice for the project. The difference between an 80:20 split and a 60:40 split is a factor of two for the smaller share and only 25% for the larger share. This observation leads towards funding splits closer to 50:50 in years when both resources are oversubscribed. Over the last several years, this split has been between 40:60 and 60:40, and for the purpose of setting project targets, the split in this range that gives the lower total performance was selected for the project targets (in case non-GPU software continued to be a sizeable fraction of the workload), thus 60% conventional and 40% GPU. In any year, however, the split would be adjusted to yield the best science for USQCD.

In FY16, the project will decommission systems purchased in 2010 and 2011, including the remaining ARRA funded systems at Jefferson Lab. There will also be some attrition of systems purchased in 2012. This reduction in capacity will include about 9 TFlops of conventional clusters and 34 TFlops of GPU capacity, a total of about 43 TFlop/s sustained. FY 2016's performance goal primarily sustains performance while allowing the mix of machines to evolve.

Sustained performance on conventional clusters is defined as the average of single precision DWF and improved staggered ("HISQ") actions on jobs utilizing 128 MPI ranks. "Linpack" or "peak" performance metrics are not considered, as lattice QCD codes uniquely stress computer systems, and their performance does not uniformly track either Linpack or peak performance metrics across different architectures. GPU clusters or other accelerated architectures are evaluated in such a way as to take into account the Amdahl's Law effect of not accelerating the full application, or of accelerating the non-inverter portion of the code by a smaller factor than the inverter, to yield an "effective" sustained teraflops, or an equivalent cluster sustained performance. Effective GPU TFlops are based on benchmarks developed in FY13 to assess the performance of the NVIDIA GPUs used on the various project clusters on HISQ, clover, and DWF applications, and reflect the clock time acceleration of entire reference applications.

On average, the project spends about 8% of the annual hardware expansion budget on file servers, and 92% on compute nodes. Using this to guide projects, the evaluations below are based up a budget of \$790K for computing hardware. Thus we are looking for a target price/performance of \$0.017 per TFlops. Total budget will be adjusted upward later to incorporate unused management reserve, and adjusted again based upon the final determination of needed disk capacity.

The goal for FY16 is to install these new resources as soon as possible, while allowing consideration of the new Intel Xeon Phi and new NVIDIA Tesla GPU. The target date for operations is thus set to Sept 30, 2016, but will be moved earlier as release dates of chips becomes clearer.

3 Hardware Options

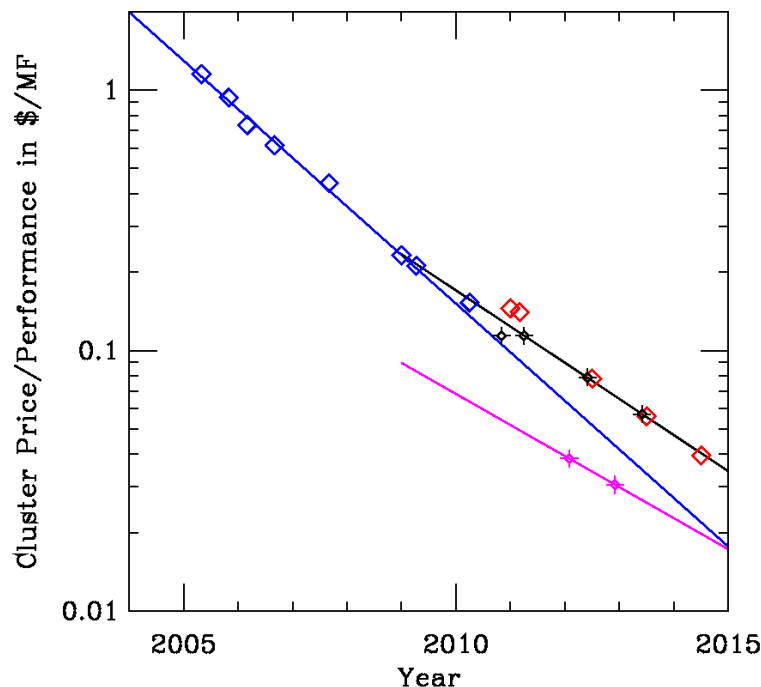
Each year the project will optimize the next procurement to yield an ensemble of hardware resources that achieves the highest performance for the portfolio of projects that USQCD intends to execute. This may include procuring two different types of computer systems in a single year.

The following types of hardware are considered in this analysis:

1. A conventional cluster, based on Intel Xeon processors with an Infiniband network.
2. A GPU accelerated cluster, based on Intel host processors, an Infiniband network, and NVIDIA GPU accelerators.
3. An Intel Xeon Phi (Knights Landing) cluster with either an Infiniband network or Intel's new Omni-Path network
4. Expansion of the current half-rack (512 nodes) IBM BG/Q system, deployed at Brookhaven National Laboratory, to a full rack (1024 nodes).

Conventional Clusters

USQCD has tracked price/performance on LQCD Infiniband-based conventional clusters deployed at Fermilab and JLab since 2005. The plot below shows these cost trends, along with exponential fits to two subsets of the data. Also included are data and an extrapolation line for GPU-accelerated clusters.



Here, the blue line is the least-squares fit to the clusters purchased between 2005 and 2011, shown as blue diamond symbols. The red diamond symbols are baseline goals used in the LQCD-ext project plan. The black line is the fit to the points from 2009 through the FY13 cluster, Bc. The magenta line connects the points corresponding to the two GPU clusters which were not memory rich, Dsg and 12k.

What is clear from this graph is that the price performance curve has a bend in 2010 such that the performance doubling time per dollar slowed from around 18 months to around 24 months. Tesla class GPUs with ECC provide roughly 4 times as much performance per dollar, but demonstrate roughly the same 24 month doubling time.

These trends continued into 2014 (the last USQCD deployed systems), but in that year more memory-rich systems were procured, which appeared to slow down the annual gain even further (a one time effect). Memory footprint has become a larger consideration as more advanced techniques yield significant performance increases in exchange for memory size, and buying more memory at the expense of flops can yield more science even if this is not (yet) reflected in our current benchmarks and metrics. It is worth noting that procurements are driven by science, and the metrics are primarily for tracking purposes.

8 and 16 GByte DIMMs are now commodity, so 64 GB/node is possible, and 128 GB/node very reasonable for conventional nodes, perhaps raising the cost of a conventional node by 10%, e.g. \$5K vs. \$4.5K. For memory rich provisioning (128 GB/node) the FY16 budget might procure ~160 nodes, for a 100% conventional system memory of 20 TB.

The largest memory footprint machine that USQCD currently operates is the Pi0 cluster at FNAL, where a 4K core job has access to 32 TB of memory (256 nodes of 128 GB, less operating system overheads). This memory size could be duplicated in FY 2016 (128 nodes of 256 GB), but the cost of doubling the memory size of Pi0 would probably be cost prohibitive on the FY16 budget. For the purposes of this analysis, then, memory size is a desirable feature, but not a driving feature.

For a high speed network, the project has for many years been choosing Infiniband fabrics, with the speed set by the performance of the node. For the last several years, QDR (quad data rate, 40 Gbps) Infiniband was the most appropriate fabric; the one exception was the 2012 quad “Kepler” GPU cluster at Jefferson Lab, where the cost of FDR (56 Gbps) was small compared to the cost per host, of order 4%, and so was selected for that small (42 node) cluster.

As x86 nodes grow in performance, higher speed fabrics can be selected, including FDR now and even 100 Gbps by next year. Alternatively, a second Infiniband card can be added to each node to support “dual rail” networking, which both doubles the bandwidth and opens up additional topologies. The cost of the second rail has been around 20% for conventional nodes, and so has not been cost effective to date.

Starting from the mid 2014 value of \$0.04/Mflops, a 2016 value of \$0.02 can be estimated. Provisioning for larger memory might raise this to \$0.025. \$790K would thus yield 31.6 TFlops. All such extrapolations are of course subject to market conditions, with errors in the tens of percent, thus 32 +/- 6 for this procurement. Clearly only using conventional nodes would not achieve our performance target.

Jefferson Lab’s procurement of Intel Haswell nodes for the experimental physics “farm” shows higher performance per core, and lower cost per core on 12 core parts and similar clock speeds. Thus peak flops continues to improve, while memory speeds improve more slowly, going from DDR3 at 1866 to DDR4 at 2133.

GPU Accelerated Clusters

For those calculations for which optimized software is available, GPU-accelerated clusters offer a substantial improvement in price/performance compared with conventional clusters. The Dsg cluster at Fermilab (delivered January 2012) and 12k at JLab (delivered November 2012) have

price/performance values of \$0.038 and \$0.031, respectively, based on a suite of application benchmarks that measure throughput with HISQ, clover, and DWF actions. By mid-FY14, using an extrapolation through these two data points, a GPU-accelerated cluster was estimated to have \$0.017/MF price/performance. The LQCD-ext II project spent more per GPU to purchase larger memory GPUs with warranties extended to 5 years at closer to \$0.028/MF on the project benchmarks. In the end the extended warranties might not be selected, so we use \$0.024 as a basis.

By the time of purchase, NVIDIA will have introduced a new generation GPU announced to have three times the performance of the K20. If a 2 year performance doubling time per dollar holds up, this would imply \$0.012/MF (and a higher cost per GPU, keeping the larger memory size). \$790K would thus yield 66 TFlops for a pure GPU cluster. Since this involves a new architecture, a higher uncertainty of 30% is used, thus 66+/-20 for the full system.

NVIDIA has announced that the Pascal GPU will have on-package high bandwidth memory (as will the new Intel Xeon Phi) and will have high bandwidth GPU-to-GPU links allowing GPUs within a node to communicate at high bandwidth without having to use the PCI bus (an important advance over the Kepler line). These anticipated improvements tend to validate the high performance projections.

For node to node communications, Infiniband remains a good choice. For such high performance nodes as quad Pascal GPUs, 100 Gbps Infiniband would be appropriate, as might multi-rail QDR or FDR Infiniband solutions. It would even be possible (but almost surely not cost effective) to install one Infiniband card per GPU, as many servers today contain 8 full bandwidth PCI slots.

The trend line for GPUs is that the integrated cost per GPU is rising (\$19K for a quad K20m node in 2012, \$24K for a quad K40 node in 2014), and so a node with a new generation of GPU might cost as much as \$30K, yielding a procurement of as few as 26 nodes, but perhaps as many as 32 nodes if prices come in low or the budget adjustments yield more available funds. If nodes are provisioned at 512 GB/node (the roughly \$2K cost of this would not be too significant), this might yield as much as 16 TB for a 100% GPU cluster.

The Pascal GPU will also have 3 point-to-point NVlinks, allowing configurations of up to 8 GPUs communicating with high bandwidth without crossing the PCI bus. Even if the configuration is held to 4 GPUs to avoid CPU bottlenecks getting worse, the much higher in-node bandwidth should allow performance per node to scale at least as well as performance per GPU. It will be necessary to benchmark real applications on this advanced in-node network to see how big an impact this will be.

Xeon Phi / Knights Landing Cluster

Unlike the previous Knights Corner generation of Xeon Phi, the Knights Landing (KNL) generation is capable of being self hosted (not an accelerator), and that is the only configuration being considered here.

Like the NVIDIA Pascal GPU, system deliveries in FY 2016 are expected.

The configuration for Knights Landing (KNL) would be single socket with “as much as” 16 GB on package memory, plus 6 DDR4 memory DIMMs. Thus 96 GB and 192 GB configurations are reasonable. If the on package memory is memory mapped and not used as a cache, this yields configurations of either 112 GB or 208 GB per node.

The first generation Xeon Phi, Knights Corner (KNC), achieved performance on “guru” level code matching the best guru code performance on a contemporary NVIDIA GPU, but due to architectural weaknesses in the KNC core did not perform as well for non-guru code. In addition, its PCI bus performance was weak. KNL is expected to address all of these weaknesses while gaining significant memory bandwidth. Core counts are publicly stated as “60+”, and if affordable parts exist at 65+ cores, software could use 64 and the O/S could have a dedicated core – this would be highly advantageous and would make reaching higher performance on linear algebra without guru code much more likely than for KNC’s 59-61 cores.

There will be 3 network options: (1) conventional PCI Infiniband HCA, (2) Intel Omni-Path PCI adapter card, and (3) integrated Intel Omni-Path. The third option might only be available with a later delivery date, but would be 100% compatible with (2).

The first option (Infiniband) would be available at 40 Gb/s, 56 Gb/s, and 100 Gb/s and would plug into the available 36 lanes of PCIe gen 3 for the non-Omni-Path-integrated KNL. As for conventional clusters, single or dual Infiniband connections could be used.

Intel Omni-Path will ship as a PCI card by Q4, 2015, x16 with 100 Gb/s in the first version (x8 available later as 58 Gb/s), including support for the OFED software stack and MPI. One port can drive 160M messages/sec, although current CPUs can’t drive that. The fabric has error detection & correction as well as QOS features that can prioritize MPI over file transfers to reduce latency. The KNL package integration will support 25 GB/s bi-directional bandwidth for each of 2 links. Port to port latency is 100-120ns with error correction. The PCI version was demonstrated at ISC’15 interconnecting Xeon and Xeon Phi nodes

Omni-Path has 48 port switches (compared with 36 ports for Infiniband), allowing 32 hosts per switch in a nominal 2:1 over-subscription configuration. This is a nice feature in that it allows for jobs of 32 sockets to run without over-subscription while leaving open the cost optimization possibility of reducing the number of switches and links for a larger cluster. 24 port switches will also be available, as will 192 and 768 port “director” switches.

In order to sell Xeon Phi, it must outperform Xeon, and a reasonable estimate would be that it does so by at least a factor of two, thus \$0.01/MFlops. Thus for \$790K, one would expect a system performance of ~80 TFlops, identical to the trend line for NVIDIA GPUs. In light of KNC’s current performance, this looks quite reasonable. Announced memory bandwidth increases to “over 400 GB/s” also makes this extrapolation credible.

Because of the possibility of running a larger class of C and C++ codes without having to re-write the software into CUDA, if the KNL chip delivers on its promise of a much improved core, a full \$790K cluster could be deployed, allowing the project to significantly exceed its target.

Software availability and maturity will be critical to moving in this direction, and all of the major USQCD collaborations are working towards this (motivated by the NERSC Cori KNL machine in 2016, as well as the ANL 2016 KNL cluster and larger ~2018 Xeon Phi / Knights Hill

capability machine). Based upon current progress, at least 2 of the 3 major code bases will have advanced solvers running by the likely time of deployment of the FY16 USQCD resource.

Initial performance on less mature software might be slightly lower than 79 TFlops, so for this document the possible system performance at launch will be conservatively assumed to be 60 TFlops, and will grow within a year to 80 TFlops as software matures.

In order to compete with the next generation dual Xeon nodes, single socket KNL nodes with twice the performance of a contemporary dual socket Xeon will need to be in the \$4K - \$5K range (depending upon node performance). If chip performance ends up much higher, cost per chip could correspondingly climb, perhaps matching the cost per GPU of \$6K - \$7K. Costs higher than that would make them non-competitive with GPUs.

As with GPUs, KNL will need to be evaluated for “effective performance”, which captures the clock time acceleration relative to a conventional node on a suite of applications, so the performance numbers referenced in this section are not pure inverter performance numbers.

The cost of doubling conventional memory from 96 GB to 192 GB is thus likely to be a 15% - 25% impact depending upon final chip cost. At the low end of cost, as many as 200 nodes of 96+16 GB could be purchased, thus 22 TB, and at the other end as few as 120 nodes of 192+16 GB, thus 25 TB. Thus Xeon Phi can achieve memory footprints similar to the Pi0 Xeon cluster if desired. (Even larger memory footprints are possible, but probably only at an additional 25% or more reduction in performance). Clearly understanding memory footprint requirements by the time of procurement award will be important.

As with x86 and GPU clusters, it will be possible to have multiple connections to the cluster interconnect. For both PCIe solutions and the integrated fabric KNL-f chip, two Omni-Path connections per socket would be feasible. Cost is unknown at this point, but assuming that it will compete with Infiniband solutions, then the second connection might be in the neighborhood of an additional 20% in cost. To date, performance improvements on conventional and GPU clusters, averaged over our portfolio of applications, have not justified a reduction in the number of nodes by 15%-20% in order to add the second network connection. As it will be for the Pascal GPU, the benefit of a second fabric connection for USQCD applications will be a part of our final system optimization.

See also: <https://software.intel.com/en-us/articles/what-disclosures-has-intel-made-about-knights-landing>, http://newsroom.intel.com/community/intel_newsroom/blog/2014/11/17/intel-reveals-details-for-future-high-performance-computing-system-building-blocks-as-momentum-builds-for-intel-xeon-phi-product.

4 Alternatives

The following sections summarize the alternative technologies considered to achieve some or all of the stated performance goals of this investment for FY16, and are listed in order of desirability.

The project team has additional NDA information beyond what is stated in this document, and there are many additional statements online, some of which may or may not be true, and these are not referenced or used in this document.

4.1 *Alternative 1: A Xeon Phi / KNL cluster released to production by Sept 30, 2016.*

Deploy and commission a KNL cluster of ~160 nodes with an initial performance of 60 Tflops, growing with software maturity to 80 Tflops, and a memory capacity of ~18 TB for a total M&S cost of \$0.79M.

The incremental lifecycle cost of this alternative is estimated as follows:

- Procure and install a 60 TF KNL cluster (\$0.79M)
- Procurement and commissioning labor, \$220K
- Operate 180 nodes using 0.125 FTE per cluster and 1 FTE / 900 nodes = 0.325 FTE, \$70K/year for a total of \$0.35M for 5 years
- Five-Year Lifecycle incremental cost: \$1.36M

Analysis: The hardware cost for this alternative is within the FY16 project budget. The potential of being able to buy an advanced architecture cluster using the full compute budget makes this the best possible alternative, both in potential total performance and in potential maximum memory footprint per job.

The price/performance of the cluster has large uncertainties as it is a new product. There is a large potential for even higher performance than stated, especially as software on this new architecture matures (which it will, in light of the large DOE deployments to come).

There is considerable space for detailed optimization of this cluster, including host memory size, network bandwidth, and network topology. As these nodes might be as much as 4 times the performance of prior x86 nodes, a 100 Gb network will almost certainly be selected, with Intel's Omni-Path being the strong favorite. The choice of memory size will be based upon more detailed data from users yet to be collected.

4.2 *Alternative 2: A 50% - 50% (by budget) mixture of conventional and GPU-accelerated clusters released to production by Sept 30, 2016*

Deploy and commission a conventional and a GPU-accelerated cluster capable of delivering respectively at least 15 TF and 39 effective TF, 54 Tflops total, at an M&S cost of \$0.79M.

The incremental lifecycle cost of this alternative is estimated as follows:

- Procure a 15 TF conventional cluster in FY16 (\$0.395M)
- Procure a 39 effective TF GPU-accelerated cluster in FY16 (\$0.395M)
- Procurement and commissioning labor, \$220K
- Operate a 16 node quad GPU cluster: $0.125 + 64/900 = 0.2$ FTE, \$42K/year, \$0.21M for five years
- Operate an 80 node conventional cluster: $.125 + 80/900 = 0.2$ FTE, \$0.21M for 5 years
- Five-Year Lifecycle incremental cost: \$1.43M

Analysis: The hardware costs for this alternative are within the FY16 project budget. The 50:50 split between conventional and GPU would be adjusted at the time of award based upon science requirements. This mixed resource would roughly replace the retiring conventional and GPU resources brought online in 2010-2011, and thus is low risk in that software is available already. Some additional improvements in software would be needed to exploit the newer on-package memory and inter-GPU links, but much of this would be expected to be contained in the QUDA package.

Because of partitioning the funds into two clusters, the largest job memory size would be much smaller, $16 \times 512 = 8$ TB or 80×128 GB = 10 TB. For some loss in performance, either of these two numbers (or both) could be doubled (16 TB or 20 TB).

4.3 Alternative 3: A pure GPU-accelerated cluster released to production by Sept 30, 2016

Deploy a GPU-accelerated cluster of up to 32 nodes sustaining 79 effective TF with a memory capacity of at least 12 TB, with an M&S cost of \$0.79M.

The incremental lifecycle cost of this alternative is estimated as follows:

- Procure a 79 effective TF GPU cluster in FY16 (\$0.79M)
- Procurement and commissioning labor, \$220K
- Operate a 32 node, 128 GPU cluster using 0.125 FTE per cluster + 1 FTE / 900 GPU = 0.267 FTE, \$57K/year, \$0.29M for 5 years
- Five-Year Lifecycle incremental cost: \$1.30M

Analysis: The hardware costs for this alternative are within the FY16 project budget. While this alternative delivers higher total performance, it yields a performance balance that is not accessible to the entire suite of USQCD supported applications, and so it might not be a good match to the collaboration's needs.

The possible memory size of the cluster ranges from 16 TB to 24 TB for 24-32 nodes with performance/footprint tradeoffs.

4.4 Alternative 4: A conventional cluster released to production by July 1, 2016

Deploy a conventional cluster by the end of June 2016 capable of sustaining at least 31 teraflop/s with an M&S cost of \$0.79M.

The incremental three-year lifecycle cost of this alternative is estimated as follows:

- Procure a 31 TFlops cluster with a memory capacity of at least 21 TB in FY16 (\$0.79M)
- Procurement and commissioning labor, \$220K
- Operate 172 nodes using $0.125 + 172/900 = 0.32$ FTE, \$69K/year, \$0.34M for five years
- Five-Year Lifecycle incremental cost: \$1.35M

Analysis: The hardware costs for this alternative are within the FY16 project budget. The overall performance of 31 TFlops falls well short of the 47 TFlops planned.

This alternative fails to meet the FY16 goal.

4.5 Alternative 5: Expand the half-rack of BG/Q deployed in Q1 2013 to a full rack, and deploy a small GPU cluster in FY16.

Expand the existing half-rack of BG/Q to a full rack in the first calendar quarter of 2016, with the expansion hardware capable of sustaining at least 16 teraflop/s for a total M&S cost of \$0.42M, and deploy a GPU cluster of 30 teraflop/s for a total M&S cost of \$0.30M, for a total of 35 teraflop/s.

The incremental three-year lifecycle cost of this alternative is estimated as follows:

- Procure parts (including spares) for a half rack BG/Q (\$0.42M)
- Procurement, installation and commissioning labor, \$200K
- Operations of the half rack BG/Q at \$70K/year for a total of \$0.35M for 5 years
- Procure a 30 effective TF GPU-accelerated cluster (\$0.30M)
- Procurement and commissioning labor, \$220K
- Operate an 12 node quad GPU cluster: \$0.16M for five years
- Five-Year Lifecycle incremental cost: \$1.77M

Analysis: The hardware costs for this alternative are within the FY16 project budget. This alternative is highly optimistic on BG/Q “fire sale” pricing, does not meet target performance, and has the largest lifecycle cost.

Because the BG/Q machine can support some applications that can’t run on GPUs, a GPU cluster is purchased to restore the balance of machines as the 2010 hardware is turned off. This alternative therefore roughly meets the hardware deployment goal 47 TFlops.

This alternative is not considered desirable as the BG/Q hardware would have to be end of life to meet this procurement budget, and even so the lifecycle costs are too high. It is included primarily as a scenario of potential interest if KNL slips by a year, or in case BG/Q parts fall in price

4.6 Alternative 6: Status Quo (no additional deployment in FY16)

Continue to operate the existing project clusters deployed at FNAL and JLab and the half rack BG/Q at BNL.

The cost of this alternative is \$1.72M in FY2016 to operate the existing facilities. The incremental cost of this alternative (new investment) is \$0.

Analysis: This alternative is included only for completeness and would not be capable of providing the necessary computational capacity to achieve the scientific goals of this project. Specifically, it would not leave USQCD with sufficient capacity to exploit the configuration generation capability of the supercomputers that DOE and NSF will have released to production during FY16.

4.7 Other Alternatives

Other alternatives may be relevant for consideration in future years. These were not considered for detailed analysis at this time, as their current state of maturity was not deemed sufficient. The alternatives include:

- Hybrid processors (CPU cores + accelerated cores): while such systems are beginning to emerge at the low end for low power devices (tablets), these are still future products for the High Performance Computing space.
- ARM64 processors: these are still too immature for serious consideration, but might evolve to be cost effective

Other mixtures of alternatives, such as part GPU and part KNL, were not considered as interesting in that segmenting a resource reduces each partition's memory capacity, makes high utilization harder to achieve, and increases system management costs. If both GPU and KNL are competitive, the project will choose the better of these two advanced architectures.

5 Discussion

The goal of this alternatives analysis is to select the purchase scenario which best optimizes the portfolio of USQCD dedicated resources. The estimates of procurement costs are only approximate; estimates of operational costs are based upon labor costs at Jefferson Lab, except for the BG/Q, which is based upon labor costs at BNL.

Budget guidance for FY16 only supports roughly replacing the lost conventional and GPU capacity, either directly or by moving to a single architecture (Xeon Phi) of similar or somewhat higher performance that has high potential for supporting a broader set of our application portfolio with only modest software investments.

Xeon Phi vs Mixed Conventional / GPU

Both of the first two alternatives (KNL and mixed conventional-GPU) have a high probability of exceeding our performance goals, and allowing real growth in USQCD resources, from the current 185 TFlops to a future ~220 TFlops, easing the pressure on a highly constrained resource.

The first, preferred, alternative has a potential for supporting a larger memory size job, and has lower operational costs. It has a higher performance uncertainty, and could thus be the alternative that yields the highest performance within budget, with performance growing as software matures.

The uncertainties in the cost per TFlops for these two alternatives leads to an approach in which the two manufacturers (NVIDIA and Intel) compete with one another on price as well as performance. The project in previous years benefitted from the Intel-AMD competition in a similar way. Furthermore, if either chip manufacturer suffers a major delay (as Intel did with the Broadwell chip), the project would be well positioned to move forward with the other solution.

Each solution has a moderately high probability of a schedule slip for silicon manufacturing, as these both involve leading edge, new processors, with on package memory. Keeping both

options open thus allows the project to quickly adapt to one or the other vendor encountering a major problem up until the point of award, weighing price/performance against delivery date slip. We have typically used 5% per month as the value of hardware today vs. the same hardware a month from today (if this were 0%, we would never spend money).

6 Conclusion

Pursue a procurement strategy that keeps both of these competitive alternatives on the table. Continue to refine the suite of benchmark applications to be used in bid evaluations, and continue to refine job memory footprint requirements for the majority of the projects' workloads. Work with silicon manufacturers to gain early access to hardware to shape the procurement as well as to evaluate performance.