

PPA Shared Cluster FY12 Strawman

FY12 FWP Hardware Strawman Plan

Intro

If the PPA FY12_FWP is approved, we will need to procure ~\$600k of computing hardware by September 30. It may not be necessary to get the hardware delivered by then, but we at least need to commit the funds.

This purchase is expected to be the first of ~5 yrs of similar scale procurements. At that point, success would be that we are at "steady state" on new /retiring hardware with floor space, power, cooling. Assuming success, the steady state system is a heterogeneous collection comprising 2-3 racks/per year plus some storage. So let's say about 20 racks total with ~3 of those in the "retiring state".

For this first year, getting the optimum mix of memory, interconnect etc. is probably the enemy of time. The 2nd year (FY13) purchase could be made in early to middle FY13, allowing for a correction.

So I propose that the first year hardware spec be the "maxed out" configuration that meets the highest needs. We can back off on memory or interconnect in FY13 to maximize core count if that is deemed useful. However, a reasonable prediction is that most use cases are headed into a multi-threaded high memory scenario, especially if we end up using virtualization.

In addition, getting as many nodes/cores as possible this first year is a good idea given the number of cores being retired soon. To help out with that, it may be possible to forgo or reduce the storage component. Specifically, KIPAC purchased (FY11) a 170TB lustre FS based on 8 OSS nodes and 2 LSI 60 disk RAID systems. The RAID boxes are expandable (up to an additional 4 at ~\$30k/ea) and the server capability can handle the additional space easily. The servers have both 10GigE and QDR Infiniband ports so they could serve both the older orange cluster clients and the new system during a transition period when both are running. So this saves a bit of both time and money, although the 10GigE infrastructure must be accounted for.

So that is the rationale.

Proposed configuration guidelines

1.) base the system on nodes that have a density of 4 cpu-sockets in 1U of space. The canonical unit here is the Dell PowerEdge C6220 Rack Server. Any density equivalent server could also be considered.

2.) The vendor has to be one that meets SLAC specs on delivery time, maintech compatability and board management compatability (eg. bios upgrades over network, IPMI interfaces, etc.).

3.) The compute nodes should be able to host both a QDR or better IB interface and a 10GigE interface. At least 2 motherboard 1GigE interfaces (or equivalent) should exist.

Note: Looks like going with 1 10GigE is cost comparable to 2x1GigE based on next gen switches being installed now.

4.) We assume a rack population of 20U (10 of the canonical units above) for planning. This leaves lots of room for rack switches etc. At the max power of ~1400W/2U this is 14kW which is near the max for Bldg 50 air cooling (I think).

Note: If water cooled racks are used can go with 15 or more units/rack -> 2 racks.

5.) Networking: Each rack could have 1 48port 1GigE for generic traffic. Assume 2-4 10GigE uplinks. Also need management switches. Also need Infiniband. Infiniband needs to be discussed. Not sure what is best there, tree or big single unit. I lean toward tree since it is more flexible for expansion in yrs 2+.

Notes: (A) Likely best choice for top-of-rack switch is 30p 10GigE with 8 uplinks, 2 per rack assuming 15 boxes.

6.) Assume E-series Intel cpus in the 8c/cpu and not top top category.

Note: Can we get prices similar to what SU got?

7.) BoE (back of envelope) cost is \$18k/2U for a canonical 4-node unit, \$500/node IB switch fabric w/cables. Top of rack is 2x30p ethernet switch per rack for 120 ports total. At \$500/port this is \$60k. IB is similar cost.

Guess cost/box is ~\$18k+4*500(eth)+4*500(IB) = \$22k.
So can get 600/22 ==> 27 boxes = 108 nodes = 1728 cores.

This appears consistent with the earlier strawman plans.

Discussion:

- Proposal is then to have SCCS (Teri Church +?) get quotes for 27 Dell C6220 servers with 4Gig/core, middle grade cpu speed etc. The nominal networking would be 10GigE using the new 7000 series switches (provided by SCCS). Also ask Teri to look at equivalent system from HP.
- Storage scenarios? Cluster storage is not long term. Where does long term data go? How paid for? Is any long term storage needed this year?
- Space,Power,Cooling planning -- Long term plan needed for ~10-15 racks, probably in one area. (What is being retired? When?)
- Use: Batch software: (LSF v8+?, condor, PBS Pro, others?, assign review team?)
- Optimizing shared use: Single thread vs. single-host vs. many host (MPI). Comparing orange cluster with generic batch config.