



Maximising job throughput using Hyper-Threading

Alastair Dewhurst, Dimitrios Zilaskos
RAL Tier1

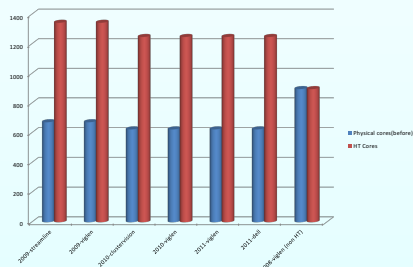
RAL Tier1 Setup

The RAL TIER1 batch farm consists of several multicore, hyperthreading capable CPUs. Increases in the amount of memory per node combined with experiences from other sites made hyperthreading an attractive option for increasing job throughput.

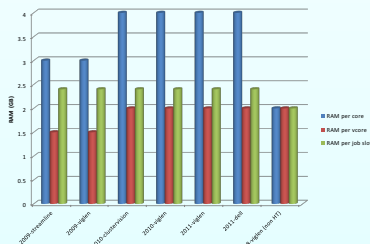
RAL supports all LHC VOs, with prime users being Atlas, CMS and LHCb, and a 10% of resources is devoted to non-LHC VOs

The virtual cores provided by hyperthreading could double the batch farm capacity, however the amount of memory available in the batch nodes did not permit that. LHC jobs require more than 2GB RAM to run smoothly.

Physical cores compared to HT cores

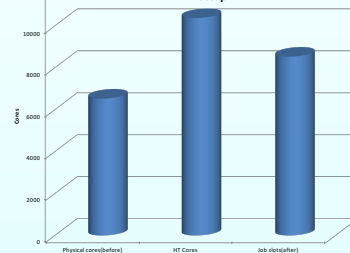


Available RAM under different setups



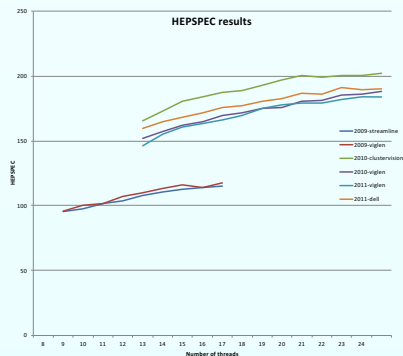
With all HT cores enabled total job slots capacity could double. In practice memory constraints resulted in an increase of 30%

Comparison of physical cores, full HT cores, and optimum setup

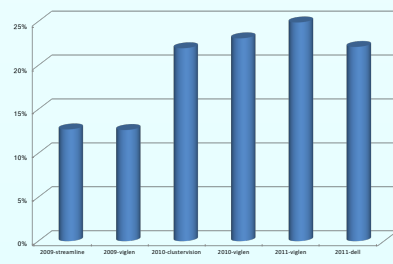


Method

- Each generation of batch farm hardware with hyperthreading capability was benchmarked with HEPSPC, progressively increasing the number of threads up to the total number of virtual cores
- Benchmarks at that time were conducted using Scientific Linux 5. Scientific Linux 6 benchmarks were run later as the batch farm was set to be updated.
- Scientific Linux 6 performed slightly better than Scientific Linux 5. The overall trend was identical
- Power, temperature, disk I/O and batch server performance were closely monitored
- The results indicated a nearly linear increase in the hepspec scores, flattening at about 14 threads for 2 CPU 4 core nodes and 20 threads for 2 CPU 6 core nodes
- The revealed sweet spots were then configured for use in the batch farm to discover where production VO jobs would perform optimally



HESPC % increase with full HT

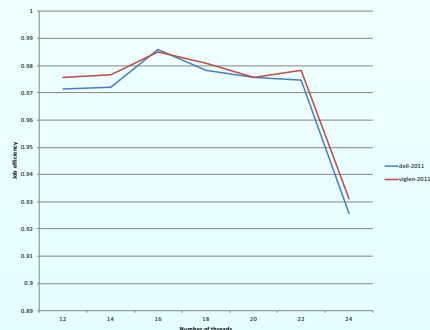


Results

- Overall, 2000 extra job slots and 9298 extra hepspec were added in the batch farm using already available hardware
- Average job time increases as expected, but overall job throughput increased
- Network/disk/power/temperature usage did not increase in a way that could negatively affect the overall throughput or require additional maneuvers
- Batch server was able to handle the extra job slots
- Of critical importance is the sharp drop in job efficiency as job slots approach the upper hyperthreading limit. This means that real world VO jobs would suffer if we went for full batchfarm HEPSPC performance!**

Make	Job Slots per WN	Efficiency	Average Job Length (mins)	Standard Deviation (mins)	Number of jobs
Dell	12	0.9715	297	376	19055
Vglen	14	0.9757	326	390	23854
Dell	14	0.9719	238	326	6118
Vglen	14	0.9767	276	341	11249
Dell	16	0.9899	345	254	6556
Vglen	16	0.9851	304	249	8716
Dell	18	0.9781	377	390	5014
Vglen	18	0.9894	356	391	6263
Dell	20	0.9758	316	346	11339
Vglen	20	0.9756	260	283	11229
Dell	22	0.9747	387	313	6317
Vglen	22	0.9783	305	236	6307
Dell	24	0.9257	544	373	6650
Vglen	24	0.9311	372	278	6713

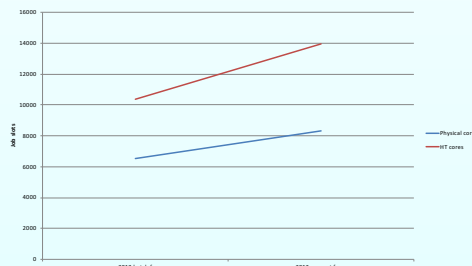
Evolution of Job efficiency as more HT cores are being used



Conclusions

- New procurements now take into account the hyperthreading capabilities
- For 2012, dual 8 core CPU systems go up to 32 virtual cores
- Systems were procured with 128 GB RAM in order to exploit full hyperthreading capabilities
- Dual Gigabit links, in the future single 10 GB as they became more cost effective
- So far RAID0 software raid setup has proven sufficient for disk I/O
- Performance gains so far on par with previous generations
- By spending a bit extra on RAM, we save more by buying fewer nodes
- This also saves machine room space, cables, and power

Evolution of batch farm size



2012 procurements benchmarks

