

# Multicore

---

Alessandra Forti

GridPP32, Pitlochry

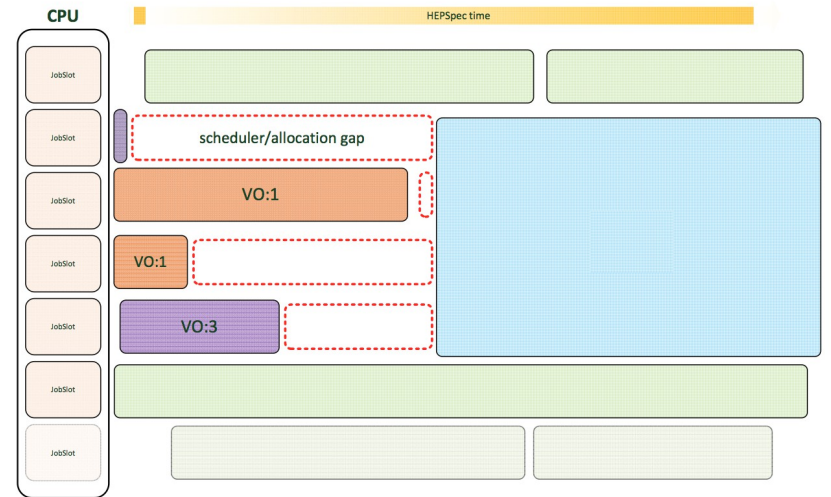
25 March 2014

# Layout

- Multicore job scheduling
- Approaches: CMS and ATLAS
- Sites
- RAL&KIT
- Backfilling
- Atlas
- Accounting
- Summary

# Multicore scheduling

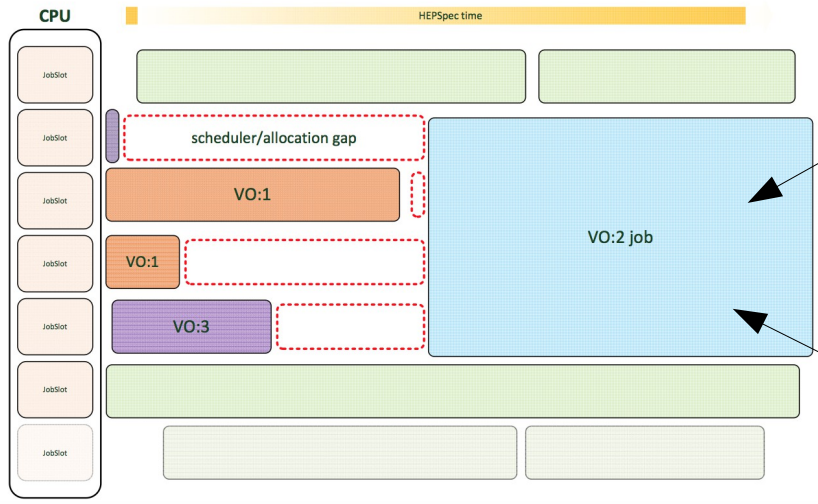
- Multicore on the WNs is “easy”
  - Jobs request all slots on 1 node
    - No need of MPI.
- The problem is running them without wasting resources.
  - The main factor in wasting resources is draining slots for multicore jobs.



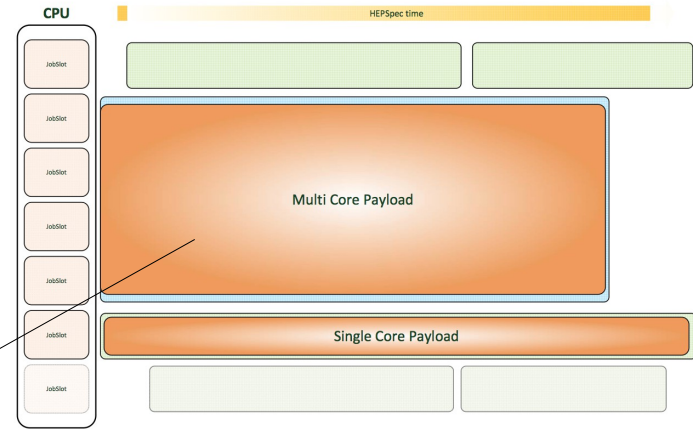
# Approaches

- 1) Scheduling multicore is **only** a site problem
  - Just request 1 node and 8 cores in the JDL
  - Blah parser will convert the request into batch system requirements
  - Batch system will handle multicore/single core jobs race conditions
  
- 2) Scheduling multicore is **also** an experiment problem
  - Pilots run different payloads until pilot walltime is exhausted
    - Dynamic scheduling (i.e. internal backfilling) done internally
  - Aim at avoiding sites allocating dedicated resources
  - Draining still present minimised by length of pilot
    - Requires quite long pilots to be effective

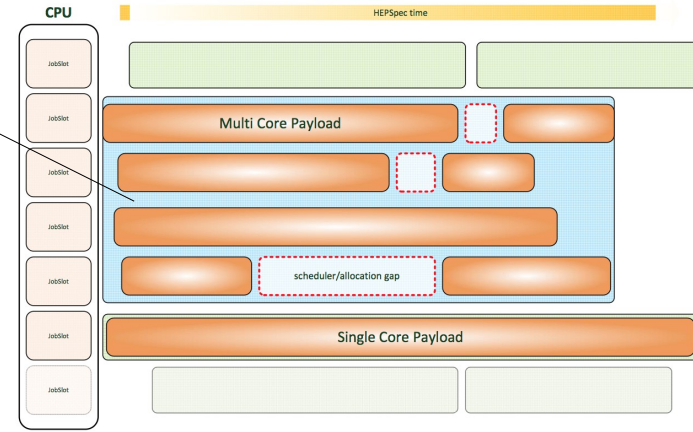
# Atlas model (1)



Inside a scheduler



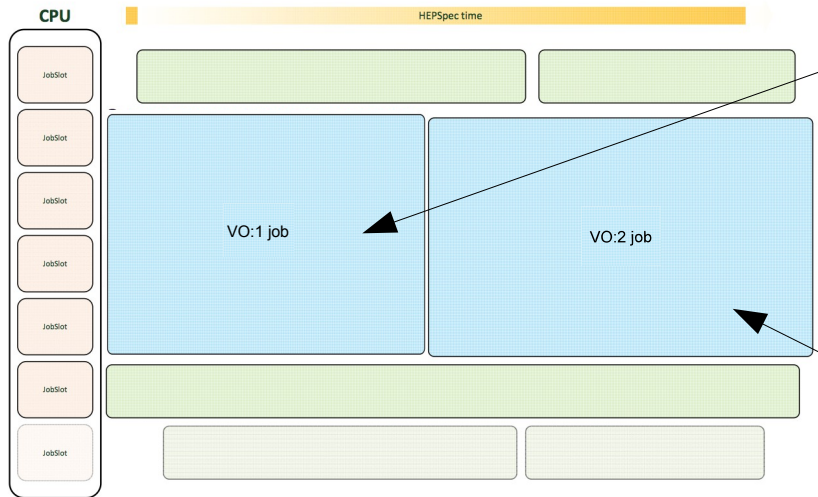
Atlas model



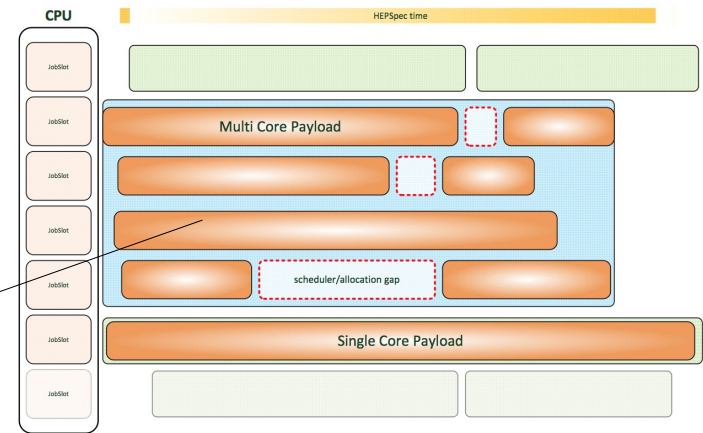
CMS model

# CMS model (2)

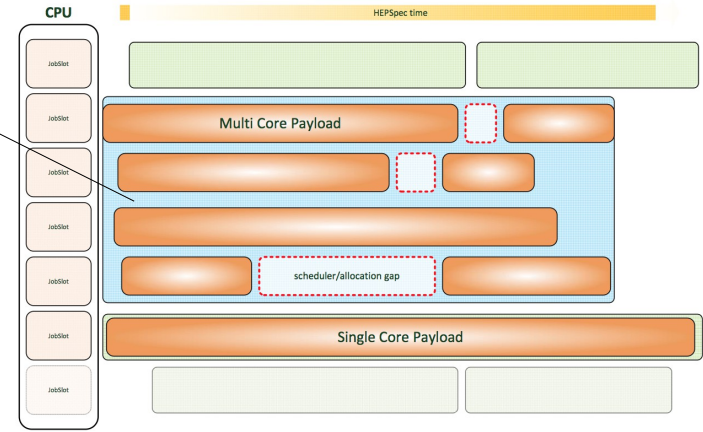
CMS model would have all Vos agree on a single pilot size



Inside a scheduler



Atlas



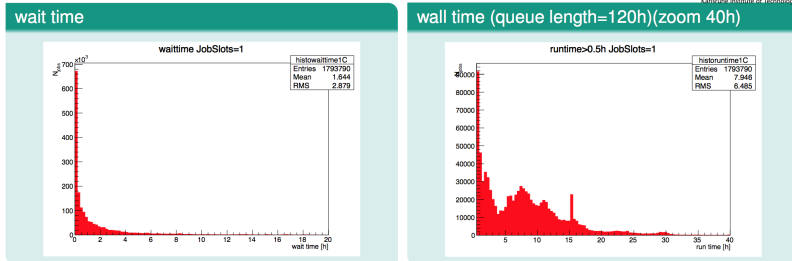
CMS

# Sites

- Atlas gave a push to run mcore before Xmas
- Several sites now accept **multicore jobs**
  - Static dedicated resources
  - Dynamic scheduling
- Different batch systems have different approaches
  - Batch system capabilities reviewed by WLCG TF
  - Torque/maui, SGE, Htcondor first review

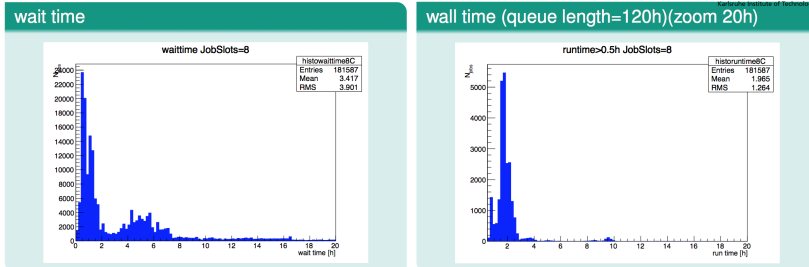
# RAL and KIT

## KIT: SingleCore Statistics 2014.Jan



Need to drain constantly

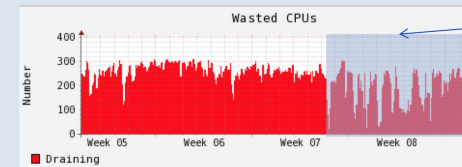
## KIT: MCore Statistics 2014.Jan



Longer wait times, shorter jobs

## Added monitoring of wasted CPUs due to draining

- Past month



Attempting to reduce wasted resources (next slide)

- We can clearly see the wastage - it's not hidden within a multi-core pilot running a mixture of single & multi-core jobs

Waste of CPU affected by submission patterns. Wavelike submission most disruptive

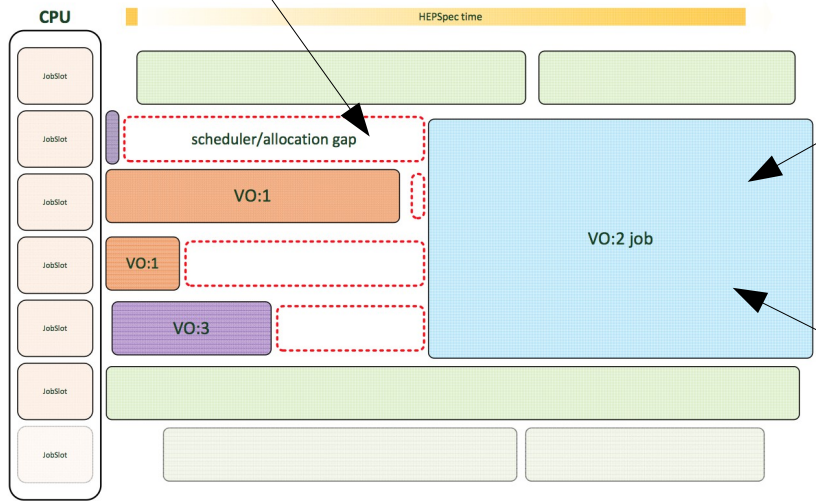
# Backfilling

- Jobs of lower priority are allowed to utilize the reserved resources only if their prospective job end (i.e. the declared wallclock usage) is before the start of the reservation
  - Most batch system are designed to do this
  - Experiments for several reasons never passed parameters to batch systems
    - One reason is also that Cream CE doesn't pass the parameters.
      - A vintage problem: since 2006!
      - Only blah SGE scripts instrumented to do it
    - ARC-CE can do it

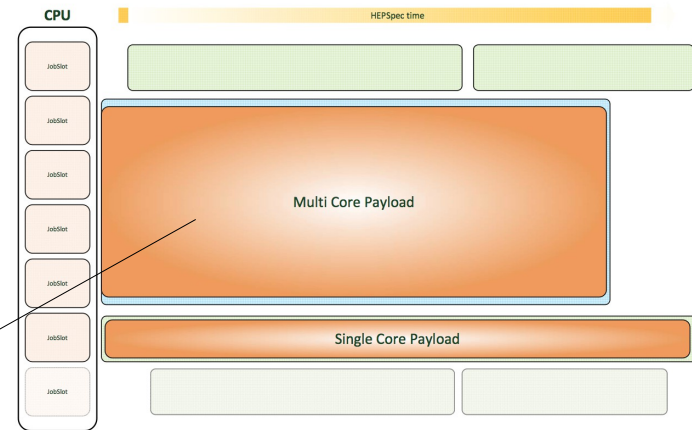
Functionality	Torque/Maui	SLURM	HTCondor	USGE/OSGE	Son of GE	LSF
Efficient Backfilling	tunable	tunable	not out-of-the-box, but similar behaviour can probably be configured	yes	yes	yes

# Most sites would like....

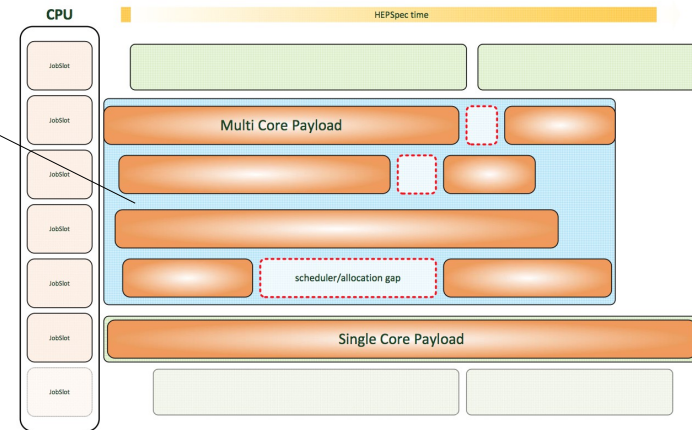
Without walltime high entropy of jobs so far was enough to fix the problem. Multicore require more organisation and backfilling with a guesstimate of walltime is needed to reduce gaps.



Inside a scheduler



Atlas model



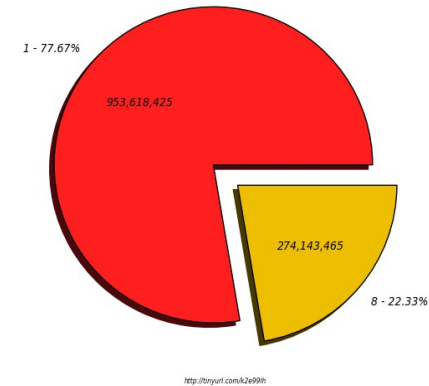
CMS model

# Atlas

- 33 sites set up
  - Efficiently or not
- Mostly at T1s and US sites
  - Large and dedicated
  - Other sites contribute too though
- Quarter of MC production done on mcore

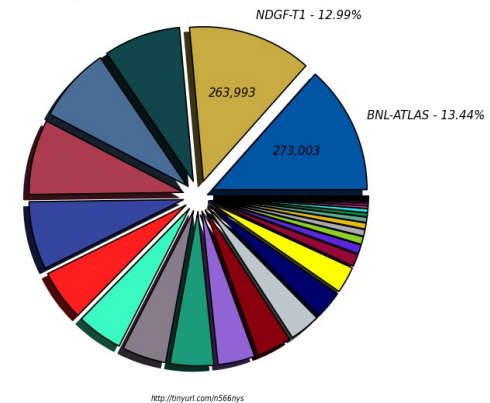
dashb@ard

WallClock HEPSPEC06 Hours (Sum: 1,227,761,891)



dashb@ard

Completed jobs Pie (Sum: 2,031,545)



BNL-ATLAS - 13.44% (273,003)	NDGF-T1 - 12.99% (263,993)
MW12 - 8.04% (163,364)	TRUMF-LCG2 - 7.88% (160,042)
LSZ-LMU - 7.87% (159,794)	BAL-CCG2 - 7.11% (144,457)
AGL12 - 5.33% (108,195)	WT2 - 4.82% (97,978)
SIQNET - 4.83% (98,037)	TRUMAN-LCG2 - 4.56% (92,581)
WUPPERTALPROD - 3.89% (78,980)	CERN-PROD - 3.63% (73,775)
UNIBE-LHEP - 3.24% (65,840)	CEKKCCG2 - 2.96% (60,138)
IN2P3-CC - 2.78% (56,474)	UKI-NORTHGRID-LANCS-HEP - 1.27% (25,854)
UKI-LTZ-OMUL - 0.94% (19,068)	NIKHEF-ELPROD - 0.79% (16,136)
UKI-FREIBURG - 0.69% (14,021)	... plus 18 more

- Multicore not only use cases to pass parameters to the batch system
  - Create a panda queue every time there is the need to pass a parameter.
    - SITE\_QUITESHORT\_ABITMOREMEM\_SOMECORES
    - ANALY\_SITE\_ABITLONGER\_LOTSOFMEM\_ALLCORES
  - Progressively messier site data structures
    - Difficult to maintain and error prone
- Task requests: walltime, memory usage, no of cores
  - Queues created dynamically
    - mcore jobs on demand on every site
    - + fast turnaround for short analysis
    - + custom jobs (high memory) can run everywhere
    - + long jobs correctly assigned to long queues

# Accounting

- Still a problem in APEL with walltime and efficiency
  - Walltime not multiplied by number of cores yet
  - Correct but would like also the walltime multiplied
- EMI-3 cream CE needed to collect correct data
  - Sites haven't upgraded yet.
- Atlas accounting has been corrected
  - Currently still testing in the prototype dashboard

# Summary

- Multicore is now ongoing
- Different experiments models
  - Aim is to find a compromise that maximise efficiency
  - Need testing with both experiments active
    - So far only Atlas
  - And shorter with walltimes declared ahead.
    - So far only BNL and Nikhef with short queues
    - Cream will be a problem for this though
- It is possible to run multicore though sites should start to play with it.
  - Some recipes already circulating