

Towards multi-parameter resource selection for HPC platforms

Master Research Thesis

Dineshkumar RAJAGOPAL
(PDES - MoSIG)

advised by

Yiannis GEORGIU

Big Data and Security(BDS) lab, BULL-SAS

September 1, 2015



1 Introduction & Motivation

2 Background

- Related Works
- SLURM Architecture
- LAYOUTS Framework

3 Resource Selection

- LAYOUTS based consumable resource selection
- Multi-parameter resource selection

4 Experiments & Performance Evaluation

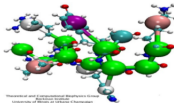
- Cons_res_layout vs Cons_res Analysis
- Multi-parameter vs Cons_res Analysis

5 Conclusion & FW

- Conclusion
- Future Work

HPC System Software Stack

- 1 Supercomputer is a **HPC Cluster**
- 2 **Usage** : Computationally intensive tasks in **Scientific Experiments** (Quantum mechanics, Weather Prediction, etc)



System Software:

Operating System, Runtime System, Resource Management, I/O System, Interfacing to External Environments

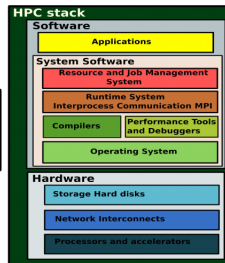
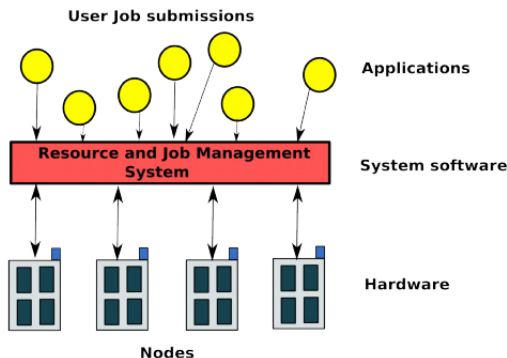


Image credit: Yiannis georgiou,BULL

What is RJMS?

- 1 The goal of a **Resource and Job Management System (RJMS)** is to satisfy users demands for computation and **assign resources to user jobs** with an efficient manner
- 2 **RJMS** knows the **complete details** about the Jobs and Resources of HPC system



Motivation of Resource Selection

- 1 **Resource selection** is an internal operation of scheduling
- 2 Due to the the evolution of HPC platform and internal nodes architecture, Resource Management is **dynamic** and **complex**
- 3 Improper resource management **hides resources information** to lose global view of resources
- 4 Power wall problem and the increasing number of nodes in HPC systems, **energy efficiency** is the important criteria
- 5 Multi-parameter resource selection to satisfy different criteria to **allocate resources perfectly**

Related Works

Feature	OAR	SLURM	FLUX
Programming Tools	Perl, MySQL	C, autoconf	C, autoconf
Purpose	Scalable & Flexible	Scalable, Flexible & Performance	Scalable, Flexible & Distributed resource & job management
Resource Management	Flat Hierarchical & managed in the Database	Linear & Managed in its own data structure(bitmap & list)	Dynamic & Flat hierarchical resource management framework
Resource Selection	Best-fit for Intel cluster architecture, SQL query to perform selection	Best-fit for different cluster architecture, custom implementation of algorithm	Next generation RJMS, only framework core functionality developed
Topology Awareness	Yes	Yes	Not Yet implemented
Internal Resource Consumption	Yes	Yes	Not Yet implemented

1 Introduction & Motivation

2 Background

- Related Works
- SLURM Architecture
- LAYOUTS Framework

3 Resource Selection

- LAYOUTS based consumable resource selection
- Multi-parameter resource selection

4 Experiments & Performance Evaluation

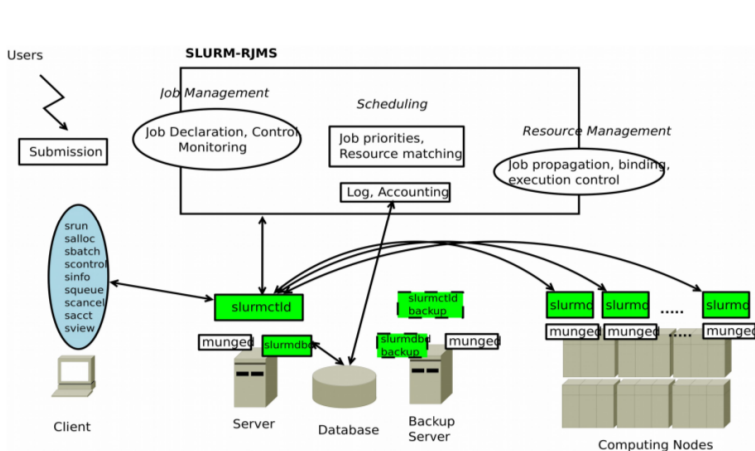
- Cons_res_layout vs Cons_res Analysis
- Multi-parameter vs Cons_res Analysis

5 Conclusion & FW

- Conclusion
- Future Work

SLURM Architecture

- 1 SLURM is an **open source RJMS** for Supercomputer

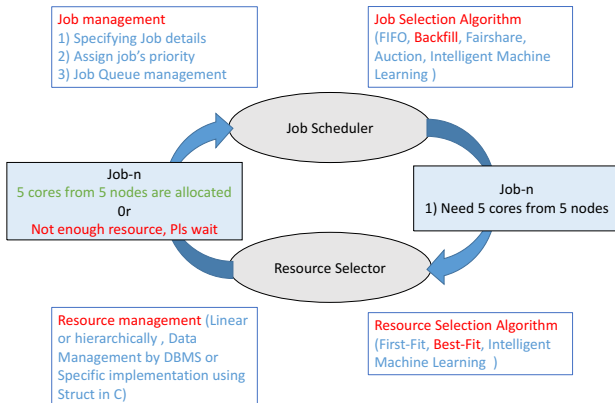


Image

credit: Yiannis georgiou,BULL

Batch Scheduling

- 1 **Scheduling behaviour** depends on the Job scheduler and Resource selector behaviour
- 2 **Resource Selection** is an internal operation of scheduling



Resource Selection Cycle

- 1 **Select** and **Topology** plugin work together to allocate **topology aware resources**
- 2 Topology plugin has information of **Switch and Node relationship in bitmap**, so comparison of two nodes list is very fast and scalable with less memory foot-print

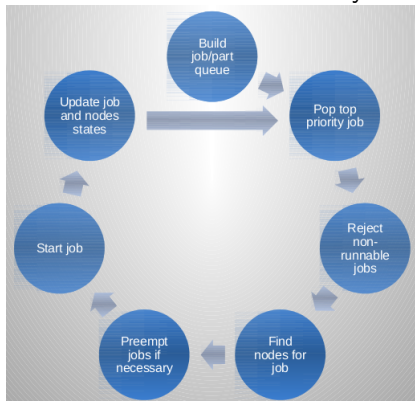


Image credit: Don Lipari, LLNL

1 Introduction & Motivation

2 Background

- Related Works
- SLURM Architecture
- LAYOUTS Framework

3 Resource Selection

- LAYOUTS based consumable resource selection
- Multi-parameter resource selection

4 Experiments & Performance Evaluation

- Cons_res_layout vs Cons_res Analysis
- Multi-parameter vs Cons_res Analysis

5 Conclusion & FW

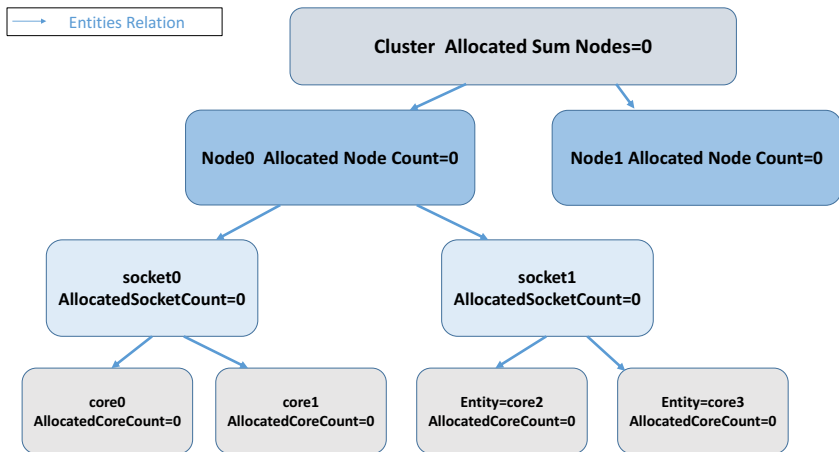
- Conclusion
- Future Work

LAYOUTS Framework

- 1 LAYOUTS is the **new resource management framework** in **SLURM v15.08**
- 2 **Any type** of entities can be manageable
- 3 Entities relation is **tree** (Inspiration of OAR resource management)
- 4 Entity attribute is called in LAYOUTS **key**, stored in the **hash key-value** format
- 5 Keeping consistency among attributes values across entities based on **keys inheritance relations**
- 6 Entities information is available in the different level of resource hierarchy to **reveal hidden information** of resources

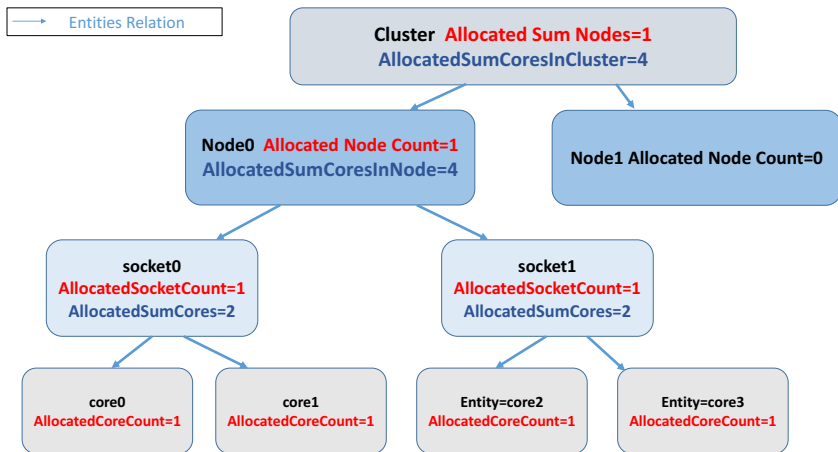
LAYOUTS Entity Keys and Key Relation

- ① If all the cores are allocated in the node, then the **node is allocated**



LAYOUTS Entity Keys and Key Relation(ctd...)

- ① If all the cores are allocated in the node, then the **node is allocated**



LAYOUTS Basic APIs

- ❶ **layouts_entity_get_kv()** - **get** key value
- ❷ **layouts_entity_set_kv()** - **set** key value
- ❸ **layouts_entity_pull_get_kv()** - **update** key relation value and **get** key value
- ❹ **layouts_entity_set_push_kv()** - **set** key value and **update** key relation

- 1 Introduction & Motivation
- 2 Background
 - Related Works
 - SLURM Architecture
 - LAYOUTS Framework
- 3 Resource Selection
 - LAYOUTS based consumable resource selection
 - Multi-parameter resource selection
- 4 Experiments & Performance Evaluation
 - Cons_res_layout vs Cons_res Analysis
 - Multi-parameter vs Cons_res Analysis
- 5 Conclusion & FW
 - Conclusion
 - Future Work

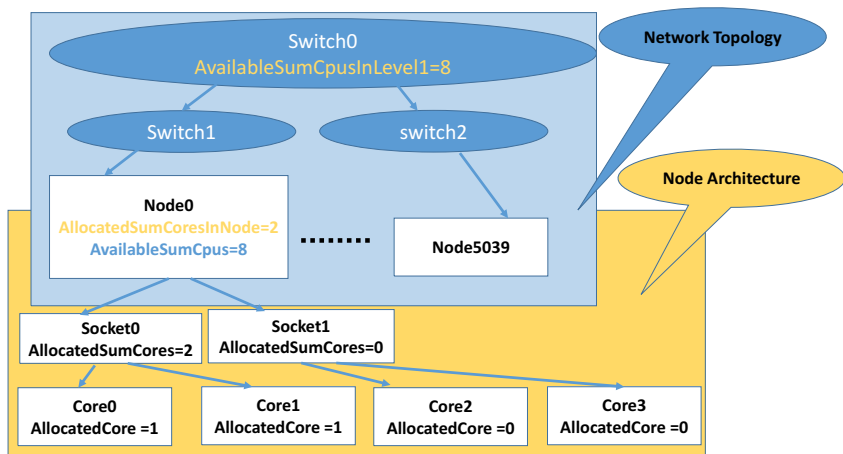
Cons_res_layout Implementation

- ❶ Cons_res is the consumable resource selection plugin
 - **Best-fit** to select **minimum satisfiable resources than maximum satisfiable resources**
 - **Topology aware** to increase the user application performance
 - **Cons_res** consumes internal resources of nodes(cores, memory, etc)
 - **Algorithm 1** in the **section 4.1** of the **report** discussed the algorithm step by step
 - **Cons_res** used **list and bitmap** to keep resource information

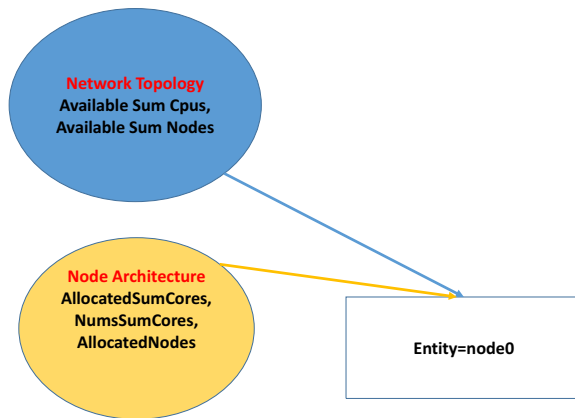
Cons_res_layout Implementation

- ❶ Cons_res is the consumable resource selection plugin
 - **Best-fit** to select **minimum satisfiable resources than maximum satisfiable resources**
 - **Topology aware** to increase the user application performance
 - **Cons_res** consumes internal resources of nodes (cores, memory, etc)
 - **Algorithm 1** in the **section 4.1** of the **report** discussed the algorithm step by step
 - **Cons_res** used **list and bitmap** to keep resource information
- ❷ **Cons_res_layout** is the **new** consumable resource selection plugin based on LAYOUTS
 - Naive implementation of **cons_res_layout** plugin performance was **25 times slower**
 - **Second version** of the **cons_res_layout** code used a combination of **bitmaps, layouts** to reach the performance of default **cons_res** code

Layouts Interrelation



Entity Global View



1 Introduction & Motivation

2 Background

- Related Works
- SLURM Architecture
- LAYOUTS Framework

3 Resource Selection

- LAYOUTS based consumable resource selection
- Multi-parameter resource selection

4 Experiments & Performance Evaluation

- Cons_res_layout vs Cons_res Analysis
- Multi-parameter vs Cons_res Analysis

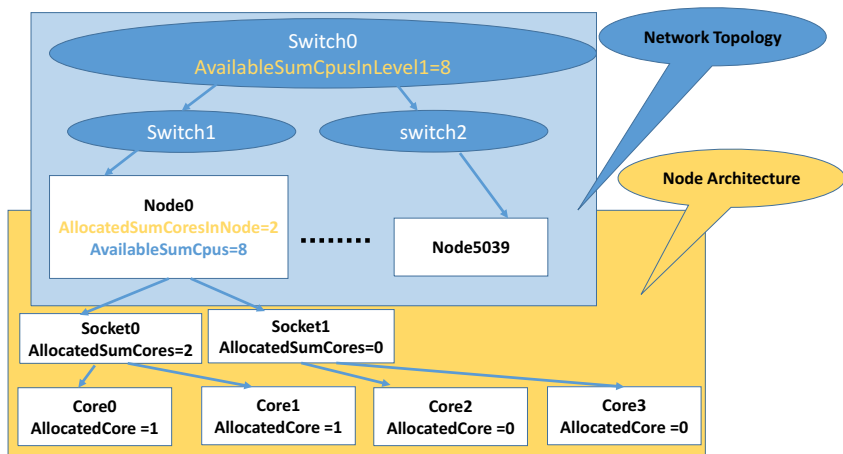
5 Conclusion & FW

- Conclusion
- Future Work

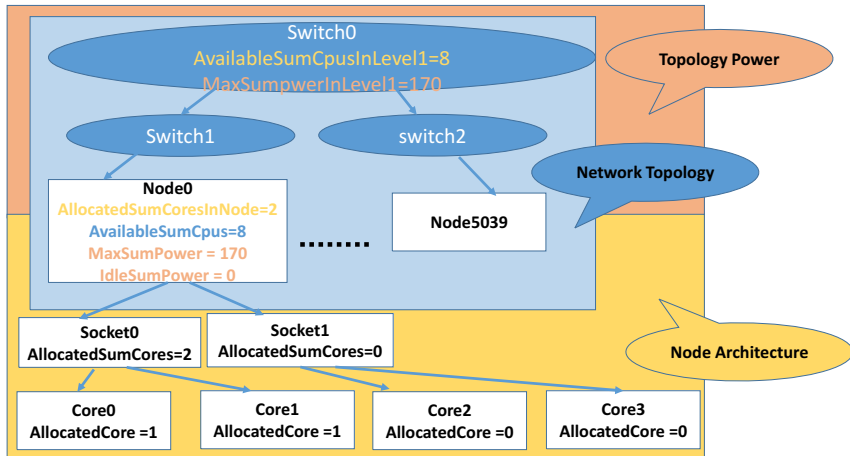
Multi-parameter Resource Selection Algorithm

- ① **Bestfit energy efficiency** to reduce energy consumption in the **Heterogeneous power consuming cluster**
 - Minimum power consuming allocated nodes
 - Already allocated nodes than idle nodes
 - Minimum power consuming idle nodes
- ② **Topology aware** to increase the user application performance
- ③ **Algorithm 2** in the **section 4.2** of the **report** discussed the algorithm step by step
- ④ **Advantage:** Multi-parameter(**Cons_res_power**) resource selection supports **user's performance** and **server's energy** criterias

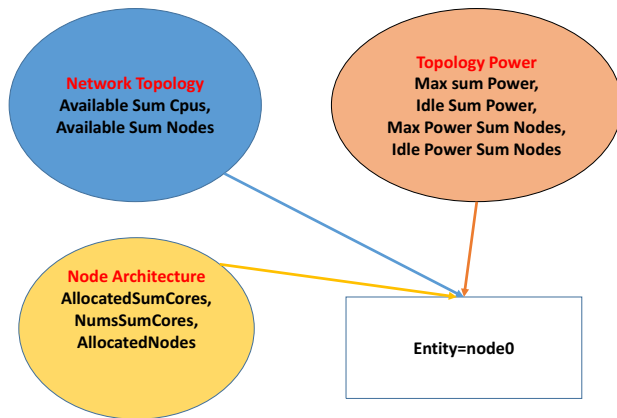
Layouts Interrelation



Layouts Interrelation(Ctd...)



Entity Global View



1 Introduction & Motivation

2 Background

- Related Works
- SLURM Architecture
- LAYOUTS Framework

3 Resource Selection

- LAYOUTS based consumable resource selection
- Multi-parameter resource selection

4 Experiments & Performance Evaluation

- Cons_res_layout vs Cons_res Analysis
- Multi-parameter vs Cons_res Analysis

5 Conclusion & FW

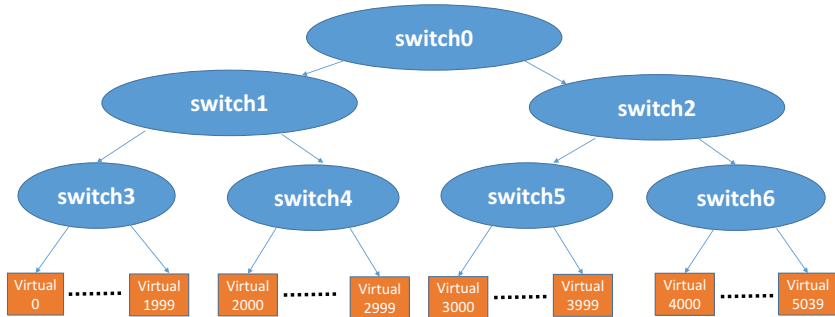
- Conclusion
- Future Work

Experiment Environment

- 1 Emulate real HPC environment using `—enable-multiple-slurmd` option
- 2 Synthetic workload of **Enhanced System Performance(ESP)** Benchmark
- 3 Use **sleep, hostname** like simple application
- 4 **Standard Workload Format(SWF)** to store workload
- 5 BULL CUZCO cluster 17 nodes to emulate **5040 nodes** HPC cluster
- 6 Each node configured as **2 sockets, 16 cores and 32GB of memory**, more details are in **Appendix C.1**

Topology Experiment Environment

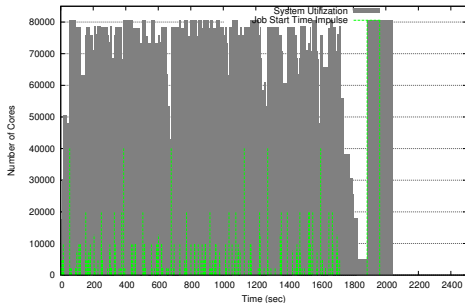
- 1 **Homogeneous** cluster environment
- 2 Simple tree topology to have **4 leaf switches** and **3 levels**



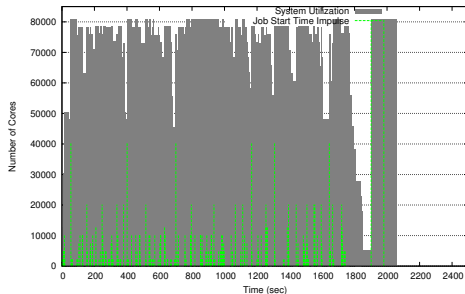
Cons_res vs Cons_res_layout system utilization

- 1 System utilization is almost same, because of following **same policy**

cons_res System utilization for Light ESP synthetic workload of 230 jobs and SLURM upon 5040 nodes (16cpu/node) cluster (emulation upon 16 physical nodes)



Layout based cons_res System utilization for Light ESP synthetic workload of 230 jobs and SLURM upon 5040 nodes (16cpu/node) cluster (emulation upon 16 physical nodes) with tree topology

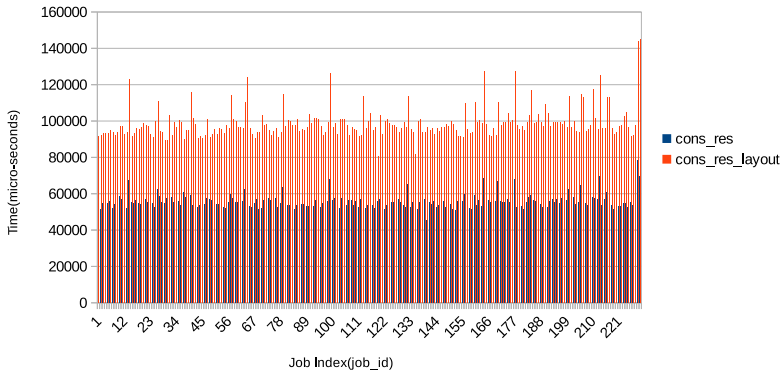


Individual Jobs Performance Comparison

- Due to large number of entity update in **Cons_res_layout** plugin, individual resource selection time increased **twice**

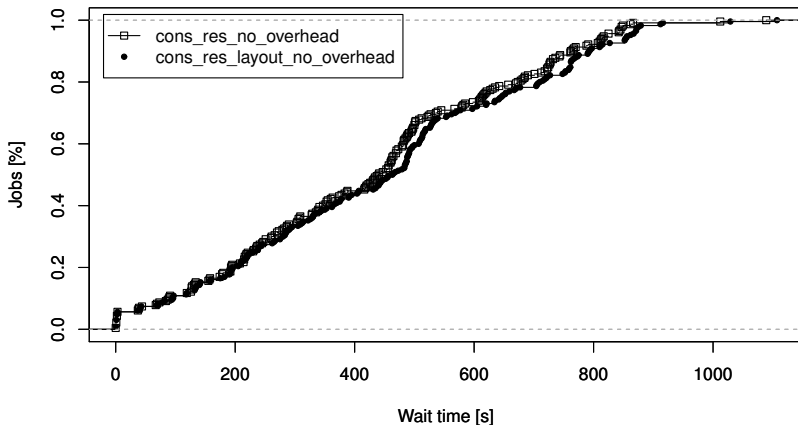
cons_res and cons_res_layout resource selection performance

(complete schedule-select cycle time in micro-seconds)



Waiting Time for 2 Plugins

**CDF on Wait time between 2 resource selection for
for Light-ESP benchmark upon a 80640 cores cluster**



1 Introduction & Motivation

2 Background

- Related Works
- SLURM Architecture
- LAYOUTS Framework

3 Resource Selection

- LAYOUTS based consumable resource selection
- Multi-parameter resource selection

4 Experiments & Performance Evaluation

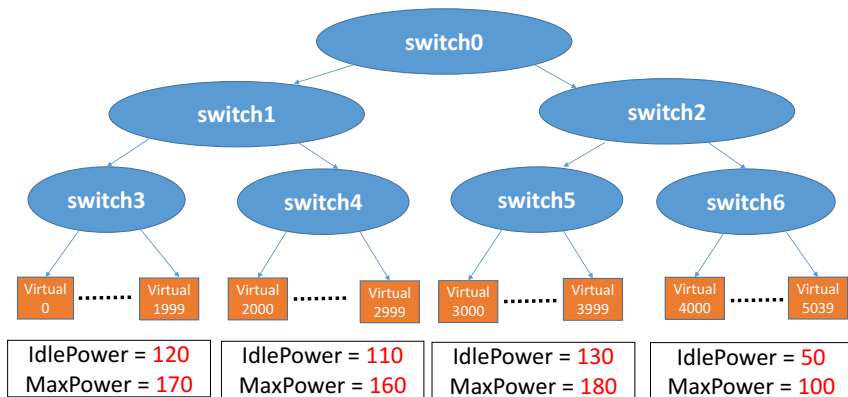
- Cons_res_layout vs Cons_res Analysis
- Multi-parameter vs Cons_res Analysis

5 Conclusion & FW

- Conclusion
- Future Work

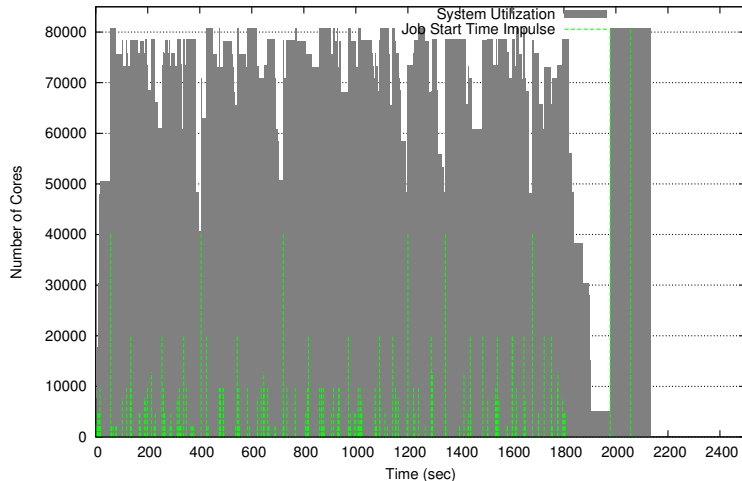
Energy Experiment Environment

- 1 **Heterogeneous power consuming nodes** cluster environment
- 2 Simple tree topology to have **4 leaf switches** and **3 levels**
- 3 **Homogeneous power consuming nodes** within leaf switches



Cons_res_power System Utilization

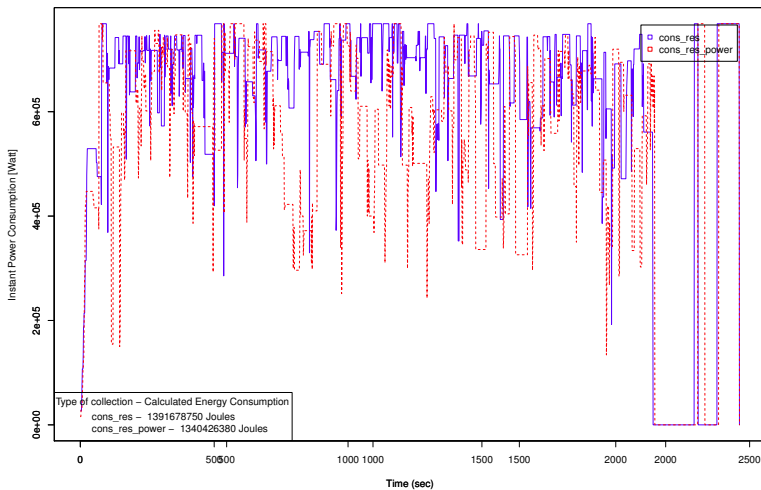
Layouts based cons_res_power System utilization for Light ESP synthetic workload of 230 jobs and SLURM upon 5040 nodes (16cpu/node) cluster (emulation upon 16 physical nodes)



Energy Efficiency

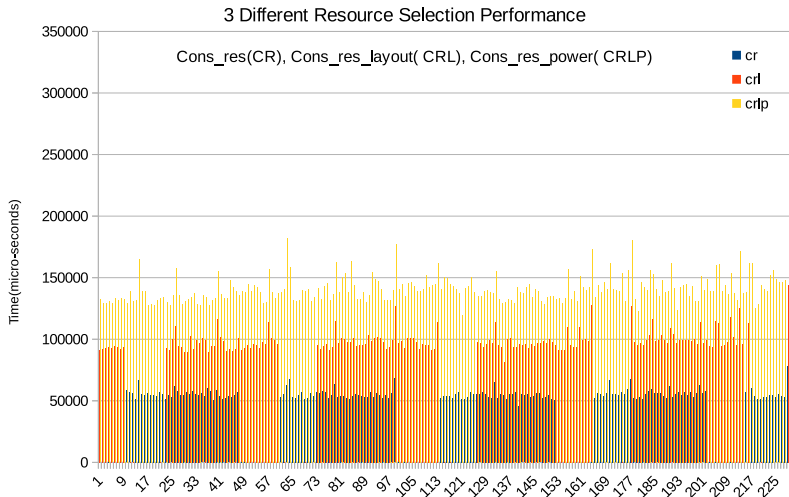
Cons_res_power energy consumption is less than Cons_res by **3.8% (51252370 Joules)** for same job workload

Power consumption comparison of 2 policies cons_res and cons_res_power



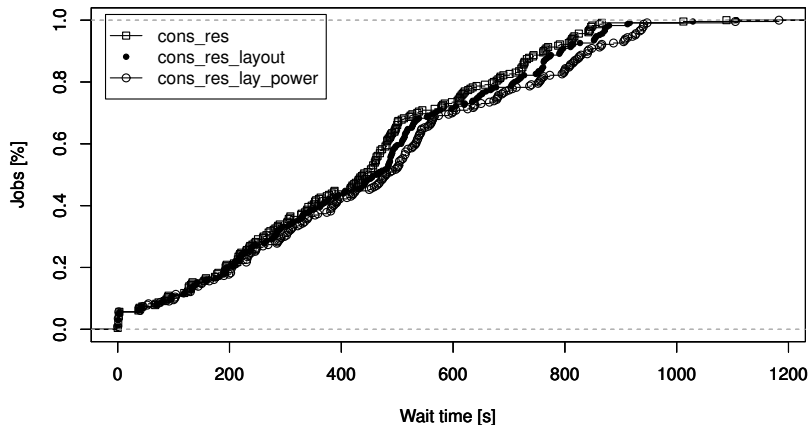
Individual Jobs Performance Comparison

- ① **Cons_res_power** individual job resource selection time increased thrice than **cons_res**



Waiting Time for 3 Plugins

CDF on Waiting time for Light-ESP benchmark upon 5040 nodes (16 cores/node) cluster comparison of 3 different cons_res



1 Introduction & Motivation

2 Background

- Related Works
- SLURM Architecture
- LAYOUTS Framework

3 Resource Selection

- LAYOUTS based consumable resource selection
- Multi-parameter resource selection

4 Experiments & Performance Evaluation

- Cons_res_layout vs Cons_res Analysis
- Multi-parameter vs Cons_res Analysis

5 Conclusion & FW

- Conclusion
- Future Work

Conclusion

- 1 SLURM **plugin enhancement** is easier using LAYOUTS resource management framework
- 2 New **Cons_res_layout** resource selection plugin developed using LAYOUTS to support **hierarchical resource management** in SLURM

Conclusion

- 1 SLURM **plugin enhancement** is easier using LAYOUTS resource management framework
- 2 New **Cons_res_layout** resource selection plugin developed using LAYOUTS to support **hierarchical resource management** in SLURM
- 3 Cons_res_layout and cons_res plugin **system utilisation and throughput** is **almost same** and **minimal increase** of resource selection performance overhead

Conclusion

- 1 SLURM **plugin enhancement** is easier using LAYOUTS resource management framework
- 2 New **Cons_res_layout** resource selection plugin developed using LAYOUTS to support **hierarchical resource management** in SLURM
- 3 Cons_res_layout and cons_res plugin **system utilisation and throughput** is **almost same** and **minimal increase** of resource selection performance overhead
- 4 Multi-parameter resource selection policy adapted in **cons_res_layout** plugin
- 5 Energy consumption of multi-parameter resource selection **reduced by 3.8 %** and **minimal increase** of resource selection performance overhead

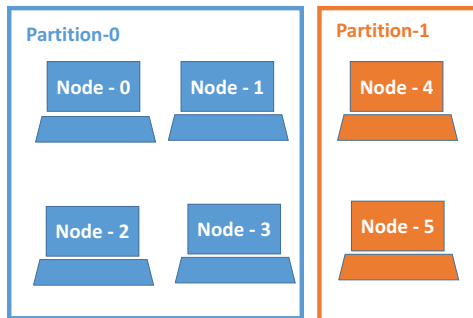
Future Work

- ① Support **partition** entities and **fat-tree** topology in the layouts plugins
- ② Adapt `cons_res_power` to support real energy values from **RAPL** or **IPMI** technique
- ③ Include **temperature** criteria in the resource selection
- ④ Experiment to measure
 - Perform **scalability experiments** with large number of nodes
 - Instantaneous job **throughput**
 - Instantaneous number of **job types** allocated
 - `Cons_res_ power` job waiting time is better than **powercapping** only approach

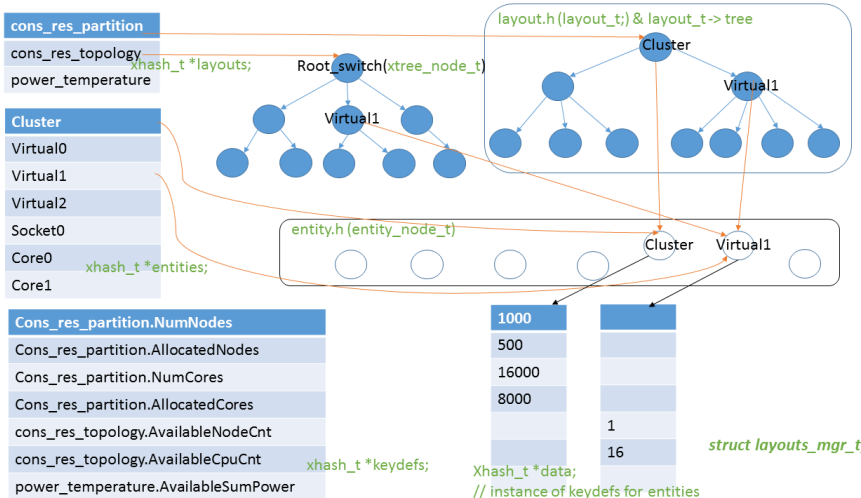
Thank you ..
Any questions?

SLURM Entities

- 1 SLURM **resource management** entities
- 2 **Job management** entities are not considered



LAYOUTS Internal Architecture



LAYOUTS Aggregate Keys

- **Child aggregate functions** - Specific operations performed on the children key value and update the calculated value in the parent's key
 - ➊ KEYSPEC_UPDATE_CHILDREN_SUM
 - ➋ KEYSPEC_UPDATE_CHILDREN_AVG
 - ➌ KEYSPEC_UPDATE_CHILDREN_MIN
 - ➍ KEYSPEC_UPDATE_CHILDREN_MAX
 - ➎ KEYSPEC_UPDATE_CHILDREN_COUNT
 - ➏ KEYSPEC_UPDATE_CHILDREN_MASK
- **Parent aggregate functions**
 - ➊ KEYSPEC_UPDATE_PARENTS_SUM
 - ➋ KEYSPEC_UPDATE_PARENTS_AVG
 - ➌ KEYSPEC_UPDATE_PARENTS_MIN
 - ➍ KEYSPEC_UPDATE_PARENTS_MAX
 - ➎ KEYSPEC_UPDATE_PARENTS_FSHARE
 - ➏ KEYSPEC_UPDATE_PARENTS_MASK

LAYOUTS New APIs

New APIs developed for **new functionality** and **performance** purpose.

- 1 **layouts_entity_get_parent_name()**
- 2 **layouts_multi_entity_set_kv()** - set **entity** same key
- 3 **layouts_multi_entity_get_kv()**
- 4 **layouts_entity_pull_get_skv()** - update **specific key** and its relations
- 5 **layouts_entity_set_push_skv()**

Cons_res_layout Performance Improvements

- ① Node and Core list is maintained in the **bitmap** data structure to compare and match another node list faster
- ② LAYOUTS to access **entities keys** in the different levels

Cons_res_layout Performance Improvements

- 1 Node and Core list is maintained in the **bitmap** data structure to compare and match another node list faster
- 2 LAYOUTS to access **entities keys** in the different levels
- 3 **layouts_multi_entity_set_kv** - transfer information from one layout to another

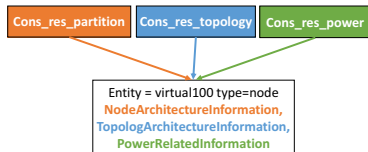
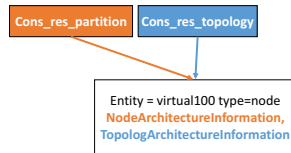
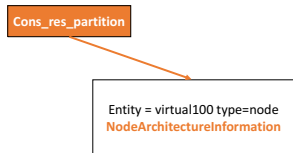
Cons_res_layout Performance Improvements

- ① Node and Core list is maintained in the **bitmap** data structure to compare and match another node list faster
- ② LAYOUTS to access **entities keys** in the different levels
- ③ **layouts_multi_entity_set_kv** - transfer information from one layout to another
- ④ **layouts_entity_pull_get_skv** - update only specific key values

Cons_res_layout Performance Improvements

- 1 Node and Core list is maintained in the **bitmap** data structure to compare and match another node list faster
- 2 LAYOUTS to access **entities keys** in the different levels
- 3 **layouts_multi_entity_set_kv** - transfer information from one layout to another
- 4 **layouts_entity_pull_get_skv** - update only specific key values
- 5 Topology aware resource selection algorithm was adapted for **normal tree** topology

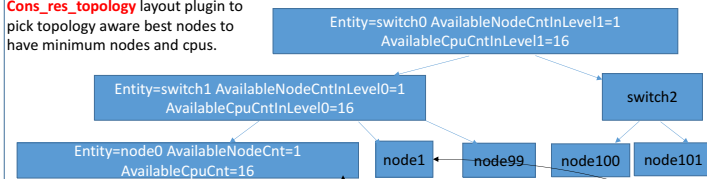
LAYOUTS Entity Data Management



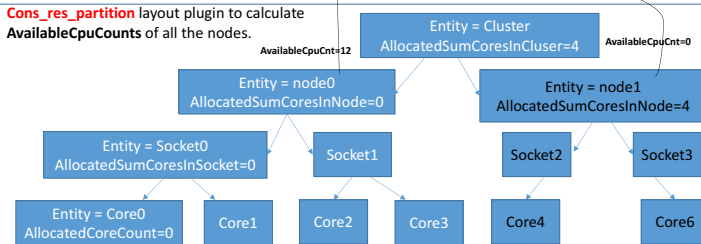
- 1 Entity and entity keys are in hash data structure and access very fast
- 2 Entity has different layouts plugin information to view entity different perspective

Cons_res_layout Architecture

Cons_res_topology layout plugin to pick topology aware best nodes to have minimum nodes and cpus.



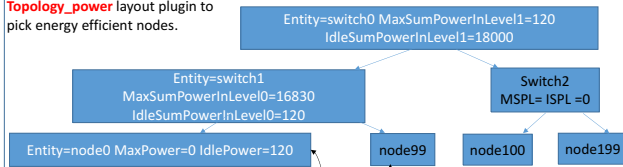
Cons_res_partition layout plugin to calculate **AvailableCpuCounts** of all the nodes.



Cons_res_layout
resource
selection plugin
to implement the
Algorithm 1.

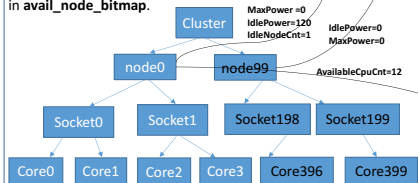
Cons_res_power Architecture

Topology_power layout plugin to pick energy efficient nodes.

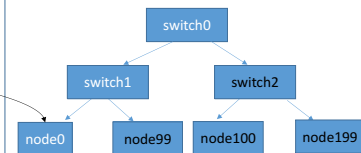


Cons_res_layout resource selection plugin enhanced to implement the **Algorithm 2**, resource selection algorithm to support energy efficient, topology aware resource selection.

Cons_res_partition layout plugin to calculate job requirement satisfying nodes **AvailableCpuCounts** and keep the selected nodes in **avail_node_bitmap**.



Cons_res_topology layout plugin to pick topology aware best nodes.



Conclusion I

- ① SLURM **plugin enhancement** is easier using LAYOUTS resource management framework

Conclusion I

- ① SLURM **plugin enhancement** is easier using LAYOUTS resource management framework
- ② **Cons_res_layout** resource selection plugin developed using LAYOUTS

Conclusion I

- ① SLURM **plugin enhancement** is easier using LAYOUTS resource management framework
- ② **Cons_res_layout** resource selection plugin developed using LAYOUTS
- ③ LAYOUTS **new APIs** developed to **increase performance** for multiple updates

Conclusion I

- ① SLURM **plugin enhancement** is easier using LAYOUTS resource management framework
- ② **Cons_res_layout** resource selection plugin developed using LAYOUTS
- ③ LAYOUTS **new APIs** developed to **increase performance** for multiple updates
- ④ Cons_res_layout and cons_res plugin **system utilisation and throughput** is **same**

Conclusion I

- 1 SLURM **plugin enhancement** is easier using LAYOUTS resource management framework
- 2 **Cons_res_layout** resource selection plugin developed using LAYOUTS
- 3 LAYOUTS **new APIs** developed to **increase performance** for multiple updates
- 4 Cons_res_layout and cons_res plugin **system utilisation and throughput** is **same**
- 5 Cons_res_layout individual resource selection performance overhead was **increased twice** than cons_res plugin
- 6 Cons_res_layout **job waiting time increased**, due to individual performance of resource selection

Conclusion II

- ① Multi-parameter resource selection policy implemented in `Cons_res_power` plugin

Conclusion II

- ① Multi-parameter resource selection policy implemented in **Cons_res_power** plugin
- ② Cons_res_power plugin achieved both **application performance** and **server energy efficiency**

Conclusion II

- ① Multi-parameter resource selection policy implemented in **Cons_res_power** plugin
- ② Cons_res_power plugin achieved both **application performance** and **server energy efficiency**
- ③ Energy consumption of multi-parameter resource selection **reduced by 3.8 %**

Conclusion II

- ① Multi-parameter resource selection policy implemented in **Cons_res_power** plugin
- ② Cons_res_power plugin achieved both **application performance** and **server energy efficiency**
- ③ Energy consumption of multi-parameter resource selection **reduced by 3.8 %**
- ④ Cons_res_power individual resource selection performance overhead was **increased thrice** than cons_res plugin
- ⑤ Cons_res_power **job waiting time increased**, due to individual performance of resource selection