

# Relation-Oriented Lattice-Model for Resource Allocation in Heterogenous Distributed Systems\*

Teng Yu<sup>†</sup>

10th Jan. 2016

---

\*M.Res. Research Project, Department of Computing, Imperial College London

<sup>†</sup>Corresponding email: *t.yu15@imperial.ac.uk*

## Abstract

This work is intended to improve the performance of resource allocation process on heterogeneous distributed systems by designing a new allocation model. It provides a novel approach to involve order theory to model the allocation process and constraints by investigating on the relations between different resource allocation requests(RAr) and mapping with the structure inside the cluster of servers, then designing ranking functions on their topology as a metric for optimisation instead of viewing the underlining problem as a multi-dimensional bin-packing instance and framing the process as a special Linear Programming formulation which is common in the literature. It focuses on the modelling process and invokes well-known algorithms in experiment to compare its performance with other models. The main contribution of this model is the ability to take resources with arbitrary allocation topology and generate a ranking function for RAr, which can then be used to create and evaluate allocation algorithms.

**Keywords:** *Resource Allocation, User Requests Relation, Modular Lattice, Birkhoff Representation, Ranking function*

## Acknowledgements

First thanks to my parents. Sincere thanks to Dr. Mark Stillwell for his countless helpful comments and advise during this project. Thanks Prof. Alexander L Wolf for his supervision and my personal tutor Dr. Sadri Fariba for advice. Thanks for the Departmental Scholarship funded by the Department of Computing, Imperial College London

Warm thanks the following friends for their helps and interesting discussions related to this project: Joel Choo (Imperial College London), Shale Xiong (Imperial College London), Jialiang Wang (Harvard University), Bichen Shi (University College Dublin), Huayi Ji (University College Cork), Bojia Ma (City University of Hong Kong) and Hongyu Yang (Eindhoven University of Technology).

# Contents

<b>1</b>	<b>Introduction</b>	<b>5</b>
<b>2</b>	<b>Background and Related Work</b>	<b>6</b>
2.1	Model on Homogenous System . . . . .	6
2.2	Model on Heterogeneous Systems . . . . .	9
2.3	Model for Experimental Environment . . . . .	12
2.4	Model from Set Theory . . . . .	13
<b>3</b>	<b>About Order and Lattice Theory</b>	<b>14</b>
3.1	Partial Orders . . . . .	14
3.2	Lattices and Valuations . . . . .	15
3.3	Birkhoff's Representation Theory . . . . .	16
<b>4</b>	<b>Basic RArS</b>	<b>16</b>
4.1	RArS Relations . . . . .	16
4.2	RArS Model . . . . .	17
<b>5</b>	<b>General RArS</b>	<b>19</b>
5.1	Non-independent RArS Model . . . . .	19
5.2	Multi-server RArS Model . . . . .	20
<b>6</b>	<b>MPC-X RArS</b>	<b>21</b>
6.1	MPC-X Topology and RAr Relations . . . . .	22
6.2	MPC-X RArS Initial Topology . . . . .	23
6.3	Birkhoff's Representation on MPC-X RAr Topology . . . . .	24
6.4	Metric on MPC-X RAr Model . . . . .	26
<b>7</b>	<b>Evaluation</b>	<b>27</b>
7.1	Evaluation Architecture . . . . .	27
7.2	Evaluation Results . . . . .	29
7.3	Future Evaluation Plan . . . . .	29
<b>8</b>	<b>Conclusions and Future work</b>	<b>30</b>
<b>9</b>	<b>Appendices</b>	<b>31</b>
9.1	VP_CPSUM Algorithms . . . . .	31
9.2	Heuristics in Wrasse Solver . . . . .	32
9.3	User Guide . . . . .	33

## List of Figures

1	System Model . . . . .	7
2	System Architecture . . . . .	12
3	2-dimensional(CPU, Memory) RAr topology . . . . .	18
4	2-dimensional(Network Link, I/O Bandwidth) RAr topology . .	20
5	Basic multi-server RAr topology . . . . .	21
6	MPC-X device . . . . .	23
7	MPC-X RAr topology . . . . .	27
8	MPC-X RAr lattice-topology . . . . .	28

# 1 Introduction

Resource Allocation is a core problem for designing high-performance heterogeneous distributed systems. From a classical point of view, most of approaches in the literature modelled this process as a multi-dimensional bin packing problem[10, 19, 21, 17, 11] whilst powerful heuristics has been developed during the last several decades[15, 5, 3].

When considering current real-world heterogeneous distributed systems, at least one new constraint should be added and one significant constraints should be removed compared with the classical case which also lead to the motivation of the work presented in this report: First, in addition of single resource allocation request (RAR) can contain multiple dimensions such as asking for some CPU capacity and memory space simultaneously, there can also be relations between RARs or say different RARs are not independent; Second, single RAR can ask for the capacity which must be provided by multiple servers simultaneously and then, this also leads to the fact that the difference between RARs may be shown by the relation between servers they needed for. Recall the traditional encoding of the resource allocation process as a typical multi-dimensional bin packing problem[19], in which we use bins to represent servers whilst balls represent services and solve the allocation as assigning balls into bins by linear programming solver. In this case, we find the first constraint we mentioned above really increase the complexity of the model as we needed either to adding more dimensions to handle with it[17, 2, 9] or update old dimensions and constraints to show the relations[21] whilst the second condition thoroughly influence the performance of the solver as not only different balls can be put in one bins, different bins can also be used for one balls.

The novel approach presented in this report shares a new light on this problem by focusing on the relations between different RARs and viewing the allocation process as a mapping between the topology of RARs and clusters of servers instead of an action like putting balls into bins. Compared with others previous approach which also involving set-variables to build up the model[18], we consider more theorems from order and lattice theory to construct the efficient model. We first define the relations between the RARs and use partially ordered set (*poset*) to represent the initial topology. We find the initial topology on basic homogenous RAR model, or say independent-RAR model is already a modular structure and can be easily ranked by a hight function while the general heterogeneous RAR model is not modular. Say if there are relations between different dimensions of RARs and multi-server RARs, the initial topology will just be a general poset which cannot be easily ranked. Then the most interesting achievement we obtained is by originally applying the Birkhoff's representation theory[7] to model this initial topology which precisely transfer the non-modular initial structure to a modular lattice through a bijection. This application actually provides the ability for our approach to model arbitrary resource allocation topology. Then, it is easy to design or invoke some ranking functions to label

the nodes in this new modular structure and obtains the metrics we needed to analysis the RArS and direct the allocation. We use the MPC-X[22] device by Maxeler Technologies as a concrete industrial example to present our result and show the generality of our theoretical analysis. we have implemented a prototype program based on our model by Java-choco solver[16] to compare the performance with the classical bin packing linear-programming model and receive satisfied result, while the main part of model implementation and real data comparison remained to be the future work of this project.

The remaining sections are described as follows: section 2 provides the detailed presentation of the background and related work of resource allocation for distributed systems and show the state-of-the-art. We illustrate the essential knowledge of order and lattice theory used in this work in section 3. The original research is described from the section 4 which gives the definitions of relations between RArS and build up the initial topology of basic independent RArS with some corresponding ranking functions. Section 5 shows the more general and complicated RArS topology which considering the new conditions and mentions the bottleneck of ranking the non-modular structure directly. Section 6 analyses our approach through a concrete industrial example, MPC-X device and illustrate the application of Birkhoff’s representation theory to overcome the bottleneck by modelling the general topology to a modular lattice and describes the ranking on it. The evaluation architecture and results are presented in section 7. We give conclusions and describe the future work in section 8.

## 2 Background and Related Work

Modelling on resource allocation problem for variant distributed systems has been widely research in the literature[19, 20, 21, 17, 9, 2, 11, 24, 25, 13, 15, 10]. In this section, we illustrate the work by Stillwell, Vivien and Casanova in 2010 [19] first as it can be viewed as a common base of the background and related work in a basic homogenous case. Then we extend the discussion to the heterogeneous case by considering servers which can contains different types of resources as illustrated in [20] and considering the time-varying and shared-resource need as shown in [17, 9, 2]. After that, we present the practical model designed in [25, 13] and briefly mention the recent work in [24] which the authors considered the work-time issue for their model. We discuss the corresponding experimental models in [5, 15] and illustrate the novel modelling approach shown in [18] which motivate our work in this report.

### 2.1 Model on Homogenous System

The work in [19] detailed illustrated a typical approach to model the problem in a subset of general distributed systems, named the *shared hosting platforms*<sup>1</sup>.

---

<sup>1</sup>Shared hosting platforms can share cluster resources among services to achieve a trade-off between high utilisation and performance isolation[23]

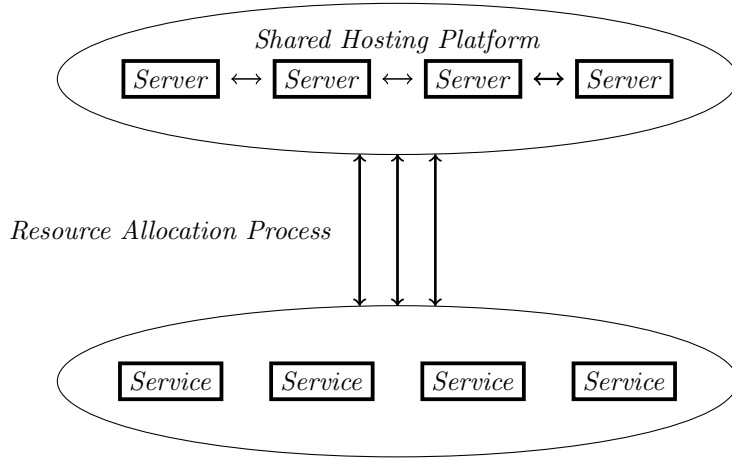


Figure 1: System Model

The authors viewed the allocation process as a multi-dimensional bin packing problem, modelled it as a Mixed Integer Linear Program then invoked different kind of algorithms to solve it and finally, evaluated the performance by some benchmark experiments. The main achievement they obtained was a formulation of this resource allocation problem in such a shared hosting platform for static workloads to make decisions when allocating hardware resources to service instances. As shown in **Figure 1**, they viewed the resource allocation as a bridge between the servers and services to construct the system model. They encapsulated homogeneous servers in cluster with high-speed switched interconnect and each server provided several resources (e.g., CPU, RAM space, I/O bandwidth, disk space) whilst they assumed each service consists of a single Virtual Machine (VM) instance<sup>2</sup>. As for resource allocation need, they classified as two types: *rigid* and *fluid*. A rigid need denoted that a specific fraction of a resource was required. The service could not benefit from a larger fraction and could not operate with a smaller fraction. A fluid need specified the maximum fraction of a resource that the service could use if alone on the server. The service could not benefit from a larger fraction, but could operate with a smaller fraction at the cost of reduced performance. For each fluid resource need, they then defined the ratio between the resource fraction allocated and the maximum resource fraction potentially used. They named this ratio the *yield* of the fluid resource need. To compare yields across services with various minimum yield requirements, The *scaled yield* of a service was as follows:

$$\text{scaled yield} = \frac{\text{yield} - \text{minimum yield}}{1 - \text{minimum yield}}.$$

<sup>2</sup>The authors also discussed an extension of their result in multi-VM services in section 6 of their paper

After that, they specified the assumptions as summarized below:

- 1) Each service consists of only a single VM instance.
- 2) Each service has constant resource needs.
- 3) The yields of all fluid resource needs are identical.
- 4) Rigid resource need are independent from fluid resource needs
- 5) User concerns, or say higher level metrics are directly related to resource fractions allocated to services

Followed that, they specified the aim precisely as *Maximize the Minimum Yield* under the two constraints below:

- 1) The resource capacities of the servers not be overcome.
- 2) A service inside a single VM instance should be allocated to a single server.

Then, the authors considered the problem formulation as a Mixed Integer Linear Program (MILP). They used  $N$  services indexed by  $i$ ,  $H$  servers indexed by  $h$  which each provided  $d$  types of resources. Fractions of these resources could be allocated to services. For each service  $i$ ,  $r_{ij}$  denoted its resource need for resource type  $j$ , as a resource fraction between 0 and 1.  $\delta_{ij}$  was a binary value that was 1 if  $r_{ij}$  was a rigid need, and 0 if  $r_{ij}$  was a fluid need. Gave  $\hat{y}_i$  to denote the minimum yield requirement of service  $i$ , a value between 0 and 1. Defined a binary variable  $e_{ih}$  that was 1 if service  $i$  run on server  $h$  and 0 otherwise. And defined  $y_{ih}$  the unscaled yield of service  $i$  on server  $h$ , which must be equal to 0 if the service did not run on the server. With these definitions the constraints of our linear program were as follows, with  $Y$  denoting the minimum yield:

$$\forall i, h \quad e_{ih} \in \{0, 1\}, \quad y_{ih} \in \mathbb{Q} \quad (1)$$

$$\forall i \quad \sum_h e_{ih} = 1 \quad (2)$$

$$\forall i, h \quad 0 \leq y_{ih} \leq e_{ih} \quad (3)$$

$$\forall i \quad \sum_h y_{ih} \geq \hat{y}_i \quad (4)$$

$$\forall h, j \quad \sum_i r_{ij}(y_{ih}(1 - \delta_{ij}) + e_{ih}\delta_{ij}) \leq 1 \quad (5)$$

$$\forall i \quad \sum_h y_{ih} \geq \hat{y}_i + Y(1 - \hat{y}_i) \quad (6)$$

This work also provided a powerful section discussing different types of algorithms they applied to solve the resource allocation problem. They analysed and evaluated the trivial solution by solving the MILP directly, 49 greedy algorithms by combining different service sorting strategies and server picking options, a genetic algorithm based on GALib library[14] and 4 vector packing algorithms<sup>3</sup>. Finally, they proved that the Choose Pack Algorithm with sorting vector lists by decreasing the sum of coordinates was the most efficient approach. We ignored the detailed algorithms analysis<sup>4</sup> and evaluation part which was far beyond the scope of this report as we focus on the modelling processing of resource allocation.

<sup>3</sup>Best Fit, First Fit, Permutation Pack and Choose Pack

<sup>4</sup>Please refer to the Appendices section for a brief illustration of the most efficient algorithm, VP\_CPSUM: Vector Packing algorithm using choose packing and sorting ordered by decreasing sum of coordinates



## 2.2 Model on Heterogeneous Systems

As a common basic case, the above work was easy to extend for heterogeneous platforms. As discussed in their following work[20], each resource provided in different servers contained two different types of capacities, elementary capacity and aggregate capacity in this heterogeneous platforms. The first capacity represented a single element in one resource dimension whilst the second denoted the total resource capacity counting all elements. For example, a server node comprised 2 cores and single memory would have different elementary and aggregate capacity for CPU but same capacity for memory; A resource allocation request might have different or same need for elementary and aggregate capacity. The aggregate capacity need might also not be an integer multiple of the corresponding elementary capacity need. Compared with the homogenous MILP model in [19], it used different vectors to represent *rigid* and *fluid* resource allocation need<sup>5</sup> instead of involving a binary indicator. It added one more constraint to handle the heterogeneous case by representing the aggregate resource capacities through the sum of resource used from all services on a server node.

While by considering other part of problems occurred in the heterogeneous systems such as shared-dimensions resource request and dynamic workload, Rai, Bhagwan, and Guha illustrated a novel approach implemented by their resource allocation solver, named *Wrasse* in [17].

Wrasse solver defined a specification language for resource allocation, which was expressive enough to encode a multitude of allocation problems without using any domain-specific abstractions. In terms of bin-packing with domain agnosticism, they assumed any ball could be assigned to any bin, balls consumed resources provided by bins and did not have different types of balls or bins map to in the problem domain as the basic case. Same as the modelling in [19, 20], Wrasse’s abstraction for resource was still a single multi-dimensional resource vector and each bin’s resource were mapped to different dimensions in that vector and encoded the constraints using a Boolean variable for each ball-bin combination then only a linear number of variables needed. In addition, by exploiting the language features that one ball could be put in only one bin, it used one integer variable for each ball and this variable was set to a value that represent its bin.

The novel approach to handle with shared resources was presented by adding new dimensions in the corresponding vectors for those resources to map to. For example, in a scenario in which there are  $n$  bins with  $k$  resources each, and  $m$  shared resources, the Wrasse resource vector has  $(k * n) + m$  dimensions. Each dimension has a capacity, fixed in the problem specification. They also gave a discussion for considering the conflicting constraints in different dimensions: say a strategy that is optimised for performance may not necessarily meet

---

<sup>5</sup>The authors renamed those two types of resource allocation needs by *requirement* and *needs* separately in this paper

fault-tolerance requirement. And trivially, if we viewed those conflicts as new shared resources constraints, we could handle them through the above way as well. A special mechanism to solve this problem more efficiently was considering *friends* and *foes* groups: they assigned related balls to the same bin and invoked user-defined function to encourage it by feedback but it cannot direct the search algorithm to explore it. Balls in a friend group are allocated to bins as few as possible while at least one ball in a foe group would be assigned to a different bin. In addition, total difference could be achieved by using pair-wise foe group. They illustrated network virtualization as a concrete example for shared resource scenario in which they extended the simple VM placement problem with network bandwidth requirements. They illustrated two previous example physical frameworks: SecondNet[9] which used the virtual data centre (VDC). Communicating VMs were placed on servers so that they did not exceed the capacities of network links connecting those servers. Each network link in the data centre was represented as resource dimension with the link capacity as the resource capacity; Oktopus[2] with Virtual Cluster(VC): a virtual star topology with  $N$  VMs: a single virtual switch connected all VMs and gave bandwidth requirement associated with each virtual link.

Next, the authors of this paper provided another novel mechanism to model the dynamic workload. Say that assigning a ball to a bin must increase the resource utilization along the appropriate resource dimensions while the static resource utilization could not model resource utilization that arised from the dynamic assignment of balls to bins. This led to the need of designing a function of the dynamic assignment. One way was also applied user-defined imperative function. There were several practical considerations when designing such a function: Soft constrains: over-constrained scenario might accept solutions where a small fraction of constrains were violated: Gave a probability with which the constraint must hold; Pinning: for ongoing process, where new balls were added, old balls evicted over time: Wrasse allowed the problem specification to 'pin' some balls to bins so as to not perturb already assigned balls in the system; Number of bins: taking the maximum of the total usage divided by total capacity across all bin-specific resources gave us a lower-bound. High-level operation: for each bin, considered all unallocated balls in parallel. A example of this case was the Microsoft Assessment and Planning Tool (MAP)[1] which took input VM time-varying requirements as a function of time. For each original resource, it created one resource dimension per time-slot. Updated the scalar resource utilization and time-varying utilization similarly. Server utilization might exceed temporally capacity by at most 10%.

The authors used real-time pair-wise bandwidth constraint as a concrete example for the case when both shared resources need and dynamic workload had to be considered. They captured the above features to design an utilization function: First stored the traffic matrix and routing information, Wrasse maintained a single variable (per dimension) for the current resource utilization; When a VM was to be placed on a server, then for every other VM on different server,

added the bandwidth requirement between this VM and the other VM on path only up to the lowest common ancestor. For other not been placed, added the requirement between this VM and the other VM on the path from this server to the root of the tree. For every VM that was placed on the same server, subtracted the bandwidth between the VMs for all links from this server to the root.<sup>6</sup> They also discussed using GPU to make best use of the corresponding hardware: Instead of using CPU, threads computed by GPU can ideally execute the same code-path (on different data), collaborate to increase sharing of the limited on-chip memory and avoid expensive synchronisation and data-dependence.

Several limitations occurred in Wrasse solver: First Wrasse cannot determine non-existence: Wrasse was sound, but not complete; Second say minimization: Wrasse did not guarantee a 'minimal' solution; Third: Wrasse's modelling was still a non-trivial transformation based on the underlining NP-complete Bin-packing problem.

Although there were still several limitations for Wrasse solver, it successfully provided a theoretical framework for modelling the general heterogeneous distributed systems. Another main problem in this field was designing the more practical framework. Hien Nguyen Van, Frederic Dang Tran, and Jean-Marc Menaud[25, 13] analysed a more complicate system model, or say architecture, to simulate the resource allocation scenario which was more closed to the really world distributed infrastructures. Instead of using bins and balls to represent, this model composed by Application Environment (AE), Local Decision Module (LDM), Global Decision Module (GDM) and datacenter which contained physical machines and VMs. Briefly, we could say that AE was faced to the application associated with specific performance goals; An application-specific LDM was associated with each AE to evaluate the process baed on the current workload using service-level metric and generated a utility function of the resource allocation; A GDM was used to interact with each LDM and the real-datacenter. It was the core of this system which contained the constraint solver to determine the management actions based on the input of LDM's utility functions and datacenter's system-level performance metrics. It worked like a black box compared with LDM. We provide an abstract system architecture graph below in **Figure 2**.

There were two utility functions in LDM: a fixed service-level function and a dynamic resource-level function which was communicated with GDM and updated for every iteration. VM allocation vectors in LDM were used for building up the upper bound constraints given by each application. As for the GDM: It involved two sequential process, one for determine VM allocation vectors for each application (VM Provisioning) and then for placing VMs and PMs and achieving optimisation (VM Packing). Beyond the constraints formulations based on

---

<sup>6</sup>For more details of the heuristics they applied in the searching process, please refer to the Appendices section in this report. The interesting approach in this part is beyond the scope of this report.

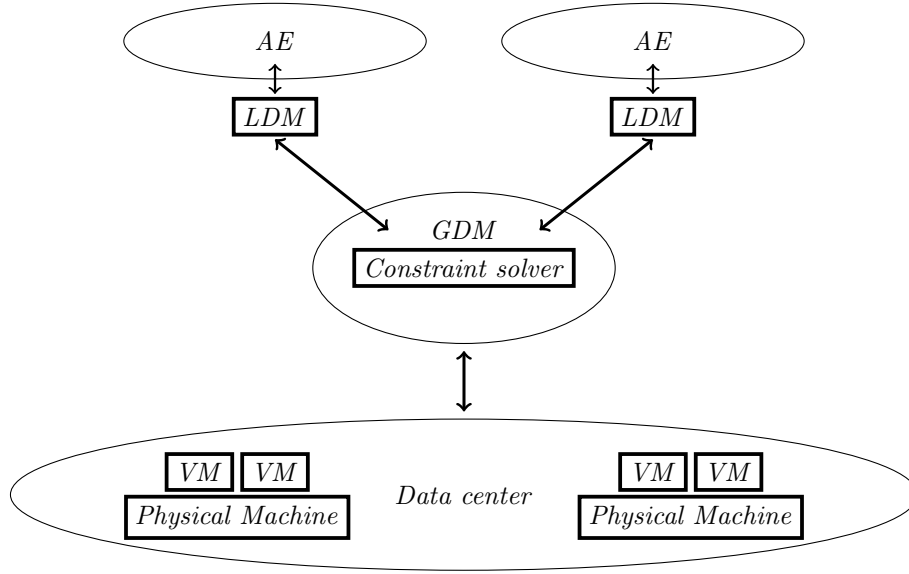


Figure 2: System Architecture

the CPU and Memory, an interesting approach in VM Provisioning was using a coefficient to allow the administrator to trade-off between the fulfilment of the performance goals and the cost of operating the required resources. Another interesting method is illustrated during the optimisation process in VM Packing: To minimise the number of migration required to reach the new VM-to-PM assignment, or say minimise the reconfiguration to provide a strategy with few interim steps and maximum degree of parallelism. In conclusion, their work presented an automatic virtual resource management system for practical distributed infrastructures. The main advantages included: it could automate the dynamic provisioning and placement of VMs; supporting for heterogeneous applications and workloads; supporting for arbitrary application topology.

By considering the working-time issue in the resource allocation problem, Shahin Vakilinia, Mustafa Mehmet Ali, and Dongyu Qiu's recent work [24] modelled the process as a job-scheduling scenario instead. It considered the randomness by involving Poisson processing during service time. Then it solved the new issues aroused by this model which including probability distribution and job blocking. It analysed the different cases for VMs, job sizes and release time, respectively to build their efficient model.

### 2.3 Model for Experimental Environment

The next important part for modelling the resource allocation were designing a efficient experimental model. Alberto Caprara and Paolo Toth's pervious work

in [5] used classes of randomly generated instances to model the 2-dimensional case. In detail, they used two variables  $c, d$  to represent the bins capacities on different dimensions, respectively. Then used  $W_j$  and  $V_j$  to represent the weights of item  $j$  on different dimensions, respectively. Finally used  $u.d.[a, b]$  to represent the uniformly distribution of the value for weights of item in the interval  $[a, b]$ . There were eight dimensions in the experimental environment. The dimensions in the first six classes were independently sampled. In classes seven and eight, the dimensions were correlated. They achieved positive correlation by setup the domain of the  $V_j$  to be a monotonic function from  $W_j$  while achieved negative correlation by setup it to be a inverse function in domain.

Rina Panigrahy, Kunal Talwar, Lincoln Uyeda, and Udi Wieder in [15] provided a more powerful approach to extend the pervious model on 2-dimensional case to multi-dimensions by setting any additional pairs of dimensions (2i-1) and (2i) correlated as dimension 1 and 2, while independent with the other dimensions. As for generating negative correlation across all dimensions: Random variables to denote the random distribution of balls in each dimension, totally 2 times of the dimensions. Multiplied with an random coefficient from  $[10, 40]$  and over two. Further noised by adding a random value from  $[0, 1]$  and ignored the overweight item.

As for defining metrics on the experimental model, Mark Stillwell in [19] sampled the resource needs from a normal probability distribution and gave the following two metrics: Failure rate ( $fr$ ): the percentage of instances for which it failed to find a solution; Distance from bound ( $dfb$ ): the difference between the achieved minimum yield and lower bound.

## 2.4 Model from Set Theory

Instead of the pervious work we mentioned above, Paul Shaw in [18] gave a new approach by modelling the underlining problem in resource allocation: the bin packing problem from a view of *Set Theory*.

The motivation of this approach was by considering a dedicated constraint and the fact that using pruning and propagation rules with lower bound method to significantly reduce search on traditional one dimensional bin packing problem. The new constraint it obtained was by considering the feature of candidate set of each bin. Trivially, there was a bound that a ball whose size was larger than the bin capacity should not be viewed as in this bin's candidate set. It gave some more precisely mathematical notations to represent the relation and generate this condition to make it dynamically during the searching process. In detail, it gives the formulas of single item elimination and commitment. A important fact it mentioned was that any bin's available space was equal to the total size to be packed, minus the loads of all other bins as it was the only rule which communicates information between different bins.

Followed this fact, it illustrated the new model of the problem from the view of *Set Theory*. As we know the goal was to assign each variable an element from its domain, then the solution could be viewed as an assignment. During the solution construction, we obtained a partial assignment which has a set of current domains at each step. Then it could view the optimisation process as a domain reduction on variables. And it clearly illustrated that its constraint included one parameter which was a vector of  $n$  element ( $n$  was the number of balls) indicating the index of the bin into which it will be placed. It also mentioned that it was NP-complete to achieve generalised arc consistency in packing constraint, so instead it described a not-complete algorithm by considering *neighbouring subsets* in candidate sets. It meant two sets were neighbouring if there was no other subset in the candidate set whose items sum to a value strictly between themselves. It implemented an linear time (even could be computed in constant time by dynamically updating) algorithm *NoSUM* to determine if there was such a interim subset exist. Then it was trivially to prune the problem instance by invoking this algorithm. The new lower bound it involved was by splitting the balls in candidate set into subsets using a constant number and applied the result to partial solutions. In conclusion, this new constraint cut search by orders of magnitude.

This approach really motivated our work in this report: our basic approach is also by modelling the relation of bins and balls by viewing one as the value in another's domain. In addition to just modelling the underlining abstract bin packing problem from the view of relations between different Sets, we try to model the actual resource allocation process as well and focus on the relations between different resource allocation requests and physical servers topology. Our approach extends the previous work by applying Order and Lattice Theory instead of just the Set theory as we model more practical constraints including non-independent requests and share-server requests by focusing on the *ordering* inside the model's topology. We first illustrate the essential knowledges used in this project in the following section.

### 3 About Order and Lattice Theory

In this section, we present the essential underlining theorems used in this work from Order and Lattice Theory. We only indicate the definitions of those theories and statements of the theorems ignoring all the verification to make the report concisely. For detail illustration of the content, we refer to Davey and Priestly's well-known book[7].

#### 3.1 Partial Orders

**Definition 1** *Let  $P$  be a set. A partial order on  $P$  is a binary relation  $\leq$  on  $P$  such that,  $\forall x, y, z \in P$ ,*  
*(i)  $x \leq x$ , (ii)  $x \leq y$  and  $y \leq x$  imply  $x = y$ , (iii)  $x \leq y$  and  $y \leq z$  imply  $x \leq z$ .*

**Definition 2** A set  $P$  equipped with an order relation  $\leq$  is said to be a partially ordered set. We use the shorthand poset in this paper.

The cover-relation between elements of a poset  $P$  is defined as follows:

$$\forall x, y \in P. x \prec y \Leftrightarrow x \leq y \text{ and } \nexists u. x \leq u \leq y.$$

We say that “ $y$  covers  $x$ ” or “ $x$  is covered by  $y$ ”. Also, the notation  $x \succeq y$  is used to indicate that  $x$  covers  $y$ , or  $x$  is equal to  $y$ .

**Definition 3** A map  $f$  from poset  $P$  onto poset  $Q$  is called an order-isomorphism iff,

$$\forall x, y \in P, x \leq y \text{ if and only if } f(x) \leq f(y) \text{ in } Q.$$

Then we say poset  $P$  and  $Q$  are isomorphic iff there exist an order-isomorphism map between  $P$  and  $Q$ .

**Definition 4** Let  $X$  be a set. The power-set  $\mathcal{Q}(X)$ , consisting of all subsets of  $X$ , is ordered by set inclusion:

$$\forall A, B \in \mathcal{Q}(X), A \leq B \text{ if and only if } A \subseteq B.$$

**Definition 5** Let  $P$  be a poset and let  $S \subseteq P$ . An element  $x \in P$  is an upper bound of  $S$  if  $s \leq x$  for all  $s \in S$ . A lower bound is defined dually. The least element of the set of all upper bounds of  $S$  is called the least upper bound and the greatest lower bound is defined dually as well.

**Definition 6** Let  $P$  be a poset and let  $S \subseteq P$ . We say  $S$  is a lower set of  $P$  iff  $\forall x \in S$  and  $y \leq x$ , then  $y \in S$ .

### 3.2 Lattices and Valuations

**Definition 7** Let  $P$  be a poset and let  $S \subseteq P$ . An element  $x \in P$  is an upper bound of  $S$  if  $s \leq x$  for all  $s \in S$ . A lower bound is defined dually. The least element of the set of all upper bounds of  $S$  is called the least upper bound and the greatest lower bound is defined dually as well.

We use  $x \sqcup y$  (read  $x$  join  $y$ ) to represent the least upper bound of  $\{x, y\}$ , whilst use  $x \sqcap y$  (read  $x$  meet  $y$ ) to represent the greatest lower bound of  $\{x, y\}$ . We use  $\bigvee S$  to represent the join of  $S$  and  $\bigwedge S$  to represent the meet of  $S$ .

**Definition 8** Let  $L$  be a non-empty ordered set.

If  $x \sqcup y$  and  $x \sqcap y$  exist for all  $x, y \in L$ , then  $L$  is called a Lattice.

**Definition 9** Given a lattice  $L$ . An element  $x \in L$  is join-irreducible iff:

$$x = a \sqcup b \text{ implies } x = a \text{ or } x = b \text{ for all } a, b \in L$$

**Definition 10** A function  $f$  on a lattice  $L$ ,  $f : L \rightarrow \mathcal{R}_0^+$  is a valuation iff

$$\forall x, y \in L. f(x \sqcap y) + f(x \sqcup y) = f(x) + f(y).$$

**Definition 11** A lattice  $L = (P, \sqsubseteq)$  is modular iff

$$\forall x, y, z \in L. x \sqsubseteq z \Rightarrow x \sqcup (y \sqcap z) = (x \sqcup y) \sqcap z.$$

The following equivalence holds in any modular lattice  $L$ :

$$(i) \forall x, y \in L. x \succ x \sqcap y \Leftrightarrow x \sqcup y \succ y.$$

$$(ii) \forall x, y, z \in L. x \succ y \Leftrightarrow x \sqcup z \succeq y \sqcup z.$$

### 3.3 Birkhoff's Representation Theory

**Definition 12** Any finite modular<sup>7</sup> lattice  $L$  is isomorphic to the lattice of lower sets of the partial order of the join-irreducible elements of  $L$ .

## 4 Basic RArS

In this section, we discuss the relation between basic homogenous resource allocation requests (RArs). We assume that RArS are independent between different dimensions in tis theoretical case and focus on the relation of RArS in each dimension to build up their foundational model.

### 4.1 RArS Relations

Consider to handle with resource that have difficult-to-represent capacities. User's request may be special on different processing element or say the relation between them. Then we find the capacity of a cluster is the set of all allocation request that it can service. Try to give a relation between RArS:

**Definition 13** Given the relation  $\leq$ , named less than, between RAr:

$$\forall \alpha \in a \text{ cluster device}, \exists A, B \in RAr, \text{ such that } A \leq B \Leftrightarrow \alpha(B) \mapsto \alpha(A).$$

We say  $\alpha(B)$  is true iff  $\alpha$  can serve B. Then in this definition, A less than B means A asks for less cluster's capacities than B. We also give the more constraint relation,  $A \preceq B$ , named A is covered by B iff A is less than B and there is no interim nodes between A and B. Then we can define a parallel relation between RArS as follow:

---

<sup>7</sup>In Birkhoff's original definition, it used *distributed lattice* instead and give a corollary saying any distributed lattice is modular. We used the followed result directly in this report as the detail mathematical description of the original theorems and corollary with the corresponding proof is far away beyond the range of this work.



**Definition 14**

$$\forall A, B \in RAr, A \not\leq B \wedge B \not\leq A \Leftrightarrow A \parallel B.$$

Followed by defining relations between user's requests, we can try to build a closed topology involving those requests. It will be more likely a lattice as we can ordered those requests by inclusion through the relation we just defined. As the topology inside a cluster device can be viewed as a lattice as well by considering that the intersection of subset servers in cluster, it leads to a possibility to achieve mappings during the allocating process.

## 4.2 RAr Model

Consider there is a common foundation of all RAr which is a non resource needed request, we represent it as  $\perp$ .

Say there are some kind of basic RAr which hold only one dimension of resource need. For instance, one RAr, named  $A_1$ , may only need some CPU capacity for computation while another, named  $A_2$ , just need some memory capacity to record information. It is easy to verify that not all cluster devices which can service  $A_1$  can service  $A_2$  and vise versa. So in this case,  $A_1$  and  $A_2$  are parallel while  $\perp \leq A_1, A_2$ . Then it is also easy to extend this instance to the real case that there will be a large number  $n$  of different dimensions of RAr, and we have the following two conditions:

$$1) \forall i, j \in n, A_i \parallel A_j$$

$$2) \forall i \in n, \perp \leq A_i$$

Now we consider the union of those RAr  $A_i$ . Some more complexed RAr may contain multi-dimensional resource allocation requests, which can be viewed as request different one-dimensional RAr simultaneously and this can be achieve by the union of  $A_i$ . In this way, we can define a n-dimensional RAr as follows:

**Definition 15** *Given  $A^n$  denote a n-dimensional RAr, then:*

$$A^n = (\forall A_i \ i \in n) \bigwedge A_i.$$

Finally, we say any kind of resource allocation request can be viewed as a n-dimensional RAr and the topology of RAr model can be built up as a partially order set through the relations we just defined. We give an instance of two-dimensional RAr topology in **Figure 3**. We use  $A_1(1)$  to represent a one-dimensional RAr which only require 1GHz CPU and  $A_2(1)$  to represent another one-dimensional RAr requiring 1GB memory. Then  $A^2(1, 1)$  is the union of them and represent a two-dimensional RAr. For easy illustration, we suppose 1GHz CPU and 1GB memory are two atom (or say entry) requests in this instance.

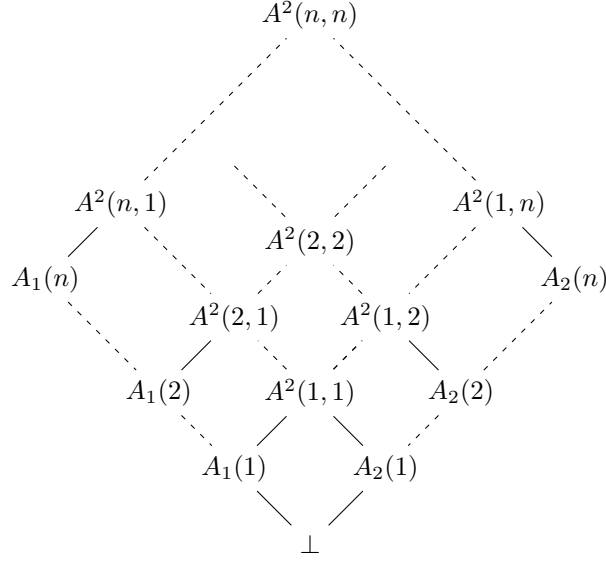


Figure 3: 2-dimensional(CPU, Memory) RArS topology

As there may be 1.5GHz CPU request in a RAr, we use dash line to denote the relation which contains intern nodes and full line to denote the closed relation.

It is easy to say the result topology is actually a modular lattice: Each pair of nodes in this topology has a least upper bound and a greatest lower bound whilst all nodes satisfy the modular law presented in Definition 11. We ignore the mathematical proof of this result.

**Note:** Instead of the theoretical analysis as above, another way of generating the topology of resource allocation requests, or more generally say *informations*, is through real data experiment. One instance of this approach can be refer to [12] which tested the web-service requests and also resulted a poset.

Then we find it is easy to give a metric on such a modular structure. Consider the value of coordinates in each node, or say vector in **Figure 3**. If we just define a ranking associated with each vector which is equal to the sum of coordinates in the vector space as follow:

**Definition 16**

$\forall A^2 \in a \text{ 2 dimensional modular structure, a ranking } R(A^2) = |A_1| + |A_2|$

Then it is obvious to say those ranking values satisfied the valuation law presented in definition 10. For example, we have  $R(A^2(1,1))=1+1=2$  and  $R(A^2(1,1))+R(\perp)=R(A_1(1))+R(A_2(1))$ . We also ignore the obvious proof. An-

other trivial ranking we can give on this structure is considering the height function:

**Definition 17**

$\forall x, y \in$  a 2 dimensional modular structure, a height  $H$  is a function that :

- (1)  $H(x) = H(y) + 1 \Leftrightarrow y \preceq x$ ;
- (2)  $H(x) = 0 \Leftrightarrow x = \perp$ .

It is not a surprise that we find the height function will give the same value for each node in this modular lattice structure as given by the ranking in definition 16. Then we say both the ranking by sum of coordinates and the height function can provide a efficient valuation for this trivial RAr's topology.

## 5 General RAr's

We begin to analysis the general case of RAr's topology in this section. In the first subsection, we extend the work in last section by reducing the constraint of independence between different dimensions in a RAr. Then in the next subsection, we consider the scenario where a RAr ask for service from multiple servers simultaneously. The generality of their RAr's topology, or say the non-modular character in their initial poset model prevent us to give a ranking on it directly as above and motivate the application of Birkhoff's representation theory to transfer the model.

### 5.1 Non-independent RAr's Model

Instead of the two independent dimensions we considered in above section, now let's focus on the case when there are two related dimensions in a single RAr. For instance, we can assume a RAr which need  $2Mb/sec$  network link's capacity and  $3Mb/sec$  I/O bandwidth. While it is easy to say any RAr containing network link's capacity need will always need at least the same I/O bandwidth. Then, we say those two dimensions in a RAr are not independent. We use  $B_1(1)$  to represent a one-dimensional RAr which only require  $1Mb/sec$  network link's capacity and  $B_2(1)$  to represent another one-dimensional RAr requiring  $1Mb/sec$  I/O bandwidth. Then  $B^2(1, 1)$  is the union of them and represent a two-dimensional RAr. Same as the above section, we suppose  $1Mb/sec$  network link's capacity and  $1Mb/sec$  I/O bandwidth are two atom requests in this instance. We formalise the relation between those two RAr's dimensions as below:

$$\forall i \in \mathcal{N}^+, B_2(i) \preceq B_1(i)$$

Then it leads to the following corollary:

$$\forall i \in \mathcal{N}^+, B_1(i) = B^2(i, i)$$

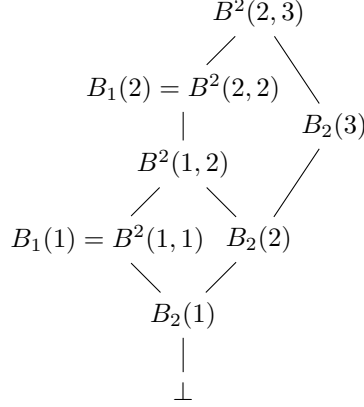


Figure 4: 2-dimensional(Network Link, I/O Bandwidth) RAr topology

Where  $\mathcal{N}^+$  denotes the set of positive integers. The basic topology of RAr with those two dimensions is shown in **Figure 4**. Compared with the topology in **Figure 3** with two independent dimensions, this topology is a really general case of a lattice instead of a modular lattice. Say if we use three nodes:  $B^2(1, 2)$ ,  $B_1(2)$  and  $B_2(3)$  as a triple, then this triple will not satisfy the modular law:  $B^2(1, 2) \preceq B_1(2)$  but  $B^2(1, 2) \sqcup (B_2(3) \sqcap B_1(2)) \neq (B^2(1, 2) \sqcup B_2(3)) \sqcap B_1(2)$ . This topology will contains more such non-modular triples or non-modular sublattice if adding the values in each dimension.

**Bottleneck of ranking the general topology:** This generality also leads to the problem that if we still want to use the rankings as mentioned in the above section, it will be failure. Give two counter examples corresponding to each of the two rankings: First consider using definition 16 as the sum of coordinates, we will have  $R(B_2(2))=2$ ,  $R(B_2(3))=3$ ,  $R(B^2(1, 2))=3$  and  $R(B^2(2, 3))=5$ , then we find  $R(B^2(1, 2))+R(B_2(3)) \neq R(B^2(2, 3))+R(B_2(2))$ . Then if we try to use the height function in definition 17, the problem is more serious that we can not even give a height value for  $B^2(2, 3)$ . Say  $H(B^2(2, 3))$  should be 5 in the direction of adding 1 from  $B^2(2, 2)$  but it should also be 4 in the direction of adding 1 from  $B_2(3)$ . In conclusion, we find the generality or say the non-modular character prevent us to rank the arbitrary allocation topology directly by some trivial ranking functions.

## 5.2 Multi-server RAr Model

We consider the case when a RAr can ask for resources provided by different servers simultaneously in this section. From a traditional bin-packing point of view, not only one bin can contains multiple balls but also one ball can be separated and put into different bins in this scenario. Instead of based on some example RAr, we analysis the relation through the general multi-server RAr

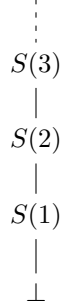


Figure 5: Basic multi-server RArS topology

directly.

It is interesting to find that when we totally refer to the relations between different multi-server RArS, the topology is much more simple than the above cases as we obtained a trivial linear-order. As shown in **Figure 5**, we use  $S(i)$  to denote a general RAr which ask for  $i$  server's resources simultaneously. The linear-order set we obtained is constructed by the trivial relation or say constraint between those RArS, say:

$$\forall i \in \mathcal{N}^+, S(i) \preceq S(i+1)$$

While the structure becomes complicate when we have to consider the relations between different servers. In the real industrial case, heterogeneous servers will provide different types or amount of resources even the number of them are the same. Further more, the difference between servers more also led by the physics topology they are allocated. For instance, two neighboured servers in a distributed system can provide resources with less network delay and enjoy more traffic tolerance than other pairs of servers. We refer to the next section as a concrete industrial example to illustrate our approach on a real-world distributed system, named MPC-X device by modelling it as a poset then transferring to a modular lattice through Birkhoff's representation theory and giving ranking on it.

## 6 MPC-X RArS

MPC-X device[22] designed by Maxeler Technologies is a novel cutting edge and non-standard computing device. MPC-X contains a cluster of Maxeler Dataflow Engines(DFE)<sup>8</sup> and each of them interconnected via a ring topology. As described in its official website[22], it provides large memory capacities (up

<sup>8</sup>As described in [6], a DFE is a physical compute resource which contains an FPGA as the computation fabric and RAM for bulk storage, and can be accessed by one or more CPU-based machines.

to 768GB) and enables remote access to DFEs by providing dual FDR/QDR<sup>9</sup> Infiniband connectivity combined with unique RDMA technology that provides direct transfers from CPU node memory to remote dataflow engines without inefficient memory copies. Other interesting features of it include the redundant power supplies and IPMI<sup>10</sup>/Lights-out management support. The computing performance of it is as powerful as a standard cluster with more than 300 cores. In detail, one task which takes 28mins on a single core of a modern machine will take only 40s on a single DFE in a MPC-X device, meaning the per-task speedup is approximately 42x. We refer to FP7 HARNESS technical paper[6] for a detailed description of MPC-X device’s architecture and performance when supporting a novel cloud computing platform.

Now we consider the MPC-X device with the additional constraints between RAr as a practical case. This is a real extended scenario that a RAr can request different servers to work for it or allocate to it simultaneously which break the classical constraints for any bin-packing problem model as they all limit that one ball can only be allocated to exactly one bin.

In this section, we first briefly show the topology of MPC-X device by a figure and then model the RAr based on its special relations. We analyse its initial RAr topology and find the bottleneck to give a ranking on this non-modular structure directly. We provide a detailed description of applying Birkhoff’s representation theory on this general poset topology to transfer to a modular lattice. Followed that, we discuss metrics on the transferred modular structure by using a ranking function and say any arbitrary allocation topology can be modelled through this approach. Finally, we analyse the time complexity of implementing the representation and constructing the ranking function. Ideally, we find our approach is within polynomial-time.

## 6.1 MPC-X Topology and RAr Relations

An abstract MPC-X device contains eight DFEs/servers<sup>11</sup> arranged in a ring as shown in **Figure 6**. We assume the servers are homogeneous in this device. Then there are two types of RAr: one for singleton servers whilst another for adjacent servers. Singleton RAr means we only need certain amount of servers to work for this request and no constraints on the physical positions of servers in the device; Adjacent RAr means we need certain amount of neighboured servers to work on it instead. Compared with the theoretical analysis in above section, this general practical scenario includes non-independent and multi-server RAr which contains all the features we discussed before.

We use  $S_n$  ( $n=1\dots 8$ ) to represent the singleton type RAr while  $A_n$  ( $n=1\dots 8$ )

---

<sup>9</sup>FDR/QDR: Four/Quad data rate

<sup>10</sup>IPMI: Intelligent Platform Management Interface.

<sup>11</sup>We use servers instead of DFEs in the following discussion to make the report more consistent.

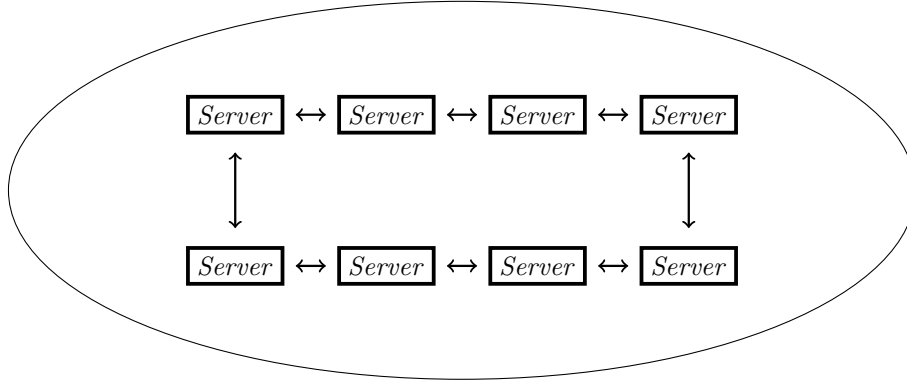


Figure 6: MPC-X device

to represent the adjacent type and in both of them, the index represent the number of servers it requested. Then it is easy to say the foundational relations below:

- 1)  $\forall i \in (1, \dots, n-1), S_i \preceq S_{i+1}$
- 2)  $\forall i \in (1, \dots, n-1), A_i \preceq A_{i+1}$
- 3)  $\forall i \in (1, \dots, n), S_i \preceq A_i$
- 4)  $A_1 = S_1$
- 5)  $A_7 = S_7$
- 6)  $A_8 = S_8$
- 7)  $A_2 \preceq S_5$
- 8)  $A_3 \preceq S_6$

## 6.2 MPC-X RArS Initial Topology

We give the initial RArS topology of the MPC-X device below **Figure 7** and ignore the totally-ordered part for easy illustration. It is totally based on the 6 classes of relations we formalised above.

**Bottleneck:** Similar with the RAr topology in **Figure 4**, it is a general lattice structure containing some non-modular substructure which prevent us to rank it directly by neither sum of vectors coordinates or a height function. For example, we cannot define the height value of  $S_5$  because two nodes,  $S_4$  and  $A_2$ , in different level are both under it immediately.

### 6.3 Birkhoff's Representation on MPC-X RAr Topology

Motivated by the bottleneck of ranking the initial non-modular RAr topology, we originally apply Birkhoff's representation theory in this scenario to map the non-modular poset topology to a downset lattice which we give the result in **Figure 8**. To illustrate the representation, we first find out all the possible downsets from the initial topology based on the definition 6 which we mentioned in section 3. For example, the downset of  $S_2$  is itself as there is no element under it while the downset of  $S_4$  is  $\{S_4, S_3, S_2\}$ . Then the easy part is to order those downsets by inclusion. Using the same example as above, we say the downset of  $S_4$  is greater than the downset of  $S_2$  as in fact, the downset of  $S_4$  include the downset of  $S_2$ . The result modular lattice topology will be constructed after we ordered all the possible downsets of the initial topology based on the above two steps.<sup>12</sup>

We use  $B_i$  to label the result graph instead of using the downset label directly to make it more concisely. To illustrate the correspondence, for example, we know  $B_1=S_2$  while  $B_{14}$  is the downset of  $S_6 \cup A_4$  which contains  $\{S_6, S_5, S_4, S_3, S_2, A_4, A_3, A_2\}$ . We provide a detail table of the representation from initial RAr topology of MPC-X device to its corresponding modular lattice topology for easy understanding in the next page. Elements in the first column are nodes in the modular lattice after Birkhoff's representation. The second column contains the name of the corresponding downset which the first column denoted. The third column represents the corresponding initial RAr to the MPC-X device which the downset in the second column contains or say the node in the first column can serve.

The time complexity of the Birkhoff's representation process on finite non-modular poset can be within polynomial time as discussed in [4]. We ignore the complexity analysis of this representation process as it is a well-known historical work in the field of discrete mathematics beyond the scope of this report.

Let's then focus on the new modular-lattice model, actually it contains some really additional *interesting* nodes which we just ignored in our first MPC-X initial RAr topology. For instance, we still use the  $B_{14}$  we just mentioned above for illustration. There cannot be any single node in the previous topology to exactly represent the scenario about the capacity need for the union of  $S_6$  and  $A_4$ . (If we say  $A_6$ , then we find that its capacity is beyond our needs as a scenario in which six servers can work whilst only five of them are adjacent can serve our need but it doesn't satisfy  $A_6$ .) On the other hand, the total sets of node in this new modular-lattice model can represent all the possible scenarios of RAr for the device based on its downward-close character.

---

<sup>12</sup>A formal description of efficient algorithms to implement Birkhoff's representation process can be refer to [4].



Node in lattice-topology	Corresponding Downset	Containing RARs
<b>B<sub>1</sub></b>	$S_2 \downarrow$	$\{S_2\}$
<b>B<sub>2</sub></b>	$S_3 \downarrow$	$\{S_3, S_2\}$
<b>B<sub>3</sub></b>	$A_2 \downarrow$	$\{A_2, S_2\}$
<b>B<sub>4</sub></b>	$(S_3 \cup A_2) \downarrow$	$\{A_2, S_3, S_2\}$
<b>B<sub>5</sub></b>	$S_4 \downarrow$	$\{S_4, S_3, S_2\}$
<b>B<sub>6</sub></b>	$(S_4 \cup A_2) \downarrow$	$\{A_2, S_4, S_3, S_2\}$
<b>B<sub>7</sub></b>	$A_3 \downarrow$	$\{A_3, A_2, S_3, S_2\}$
<b>B<sub>8</sub></b>	$(S_4 \cup A_3) \downarrow$	$\{A_3, A_2, S_4, S_3, S_2\}$
<b>B<sub>9</sub></b>	$S_5 \downarrow$	$\{A_2, S_5, S_4, S_3, S_2\}$
<b>B<sub>10</sub></b>	$(S_5 \cup A_3) \downarrow$	$\{A_3, A_2, S_5, S_4, S_3, S_2\}$
<b>B<sub>11</sub></b>	$A_4 \downarrow$	$\{A_4, A_3, A_2, S_4, S_3, S_2\}$
<b>B<sub>12</sub></b>	$(S_5 \cup A_4) \downarrow$	$\{A_4, A_3, A_2, S_5, S_4, S_3, S_2\}$
<b>B<sub>13</sub></b>	$S_6 \downarrow$	$\{A_3, A_2, S_6, S_5, S_4, S_3, S_2\}$
<b>B<sub>14</sub></b>	$(S_6 \cup A_4) \downarrow$	$\{A_4, A_3, A_2, S_6, S_5, S_4, S_3, S_2\}$
<b>B<sub>15</sub></b>	$A_5 \downarrow$	$\{A_5, A_4, A_3, A_2, S_5, S_4, S_3, S_2\}$
<b>B<sub>16</sub></b>	$(S_6 \cup A_5) \downarrow$	$\{A_5, A_4, A_3, A_2, S_6, S_5, S_4, S_3, S_2\}$
<b>B<sub>17</sub></b>	$A_6 \downarrow$	$\{A_6, A_5, A_4, A_3, A_2, S_6, S_5, S_4, S_3, S_2\}$

**Note:** A general result we can infer from this instance is that the RAr topology from any practical or say industrial distributed systems can be modelled as a modular-lattice by applying Birkhoff's representation theorem.

## 6.4 Metric on MPC-X RAr Model

At this step, we can begin to consider a metric on this modular lattice topology to represent the distance between different RAr to direct the allocation. A common metric on modular-lattice is a *valuation*. Recall the formula of valuation in definition 10. A function, or say labelling method which satisfying the valuation law will efficiently show the distance between different RAr as the sum of two different RAr's label will equal to the sum their least upper bound and greatest lower bound. In another word, valuation on the nodes of RAr's lattice topology keep the consistency of ability that to what level of RAr it can serve as shown in the corresponding downset.

A easy method to compute a valuation is by a ranking function. This function can be viewed as a formal edition of the height function we defined in definition 17 as actually they are the same. We refer to Gratzner's book [8] for the proof that the following ranking function satisfy the valuation law as it is beyond the range of this work.

**Definition 18** *Let  $L$  be a modular lattice. The ranking  $f$  on  $L$  is a function  $f : L \rightarrow \mathcal{R}_0^+$  that:  $\forall x, y \in L$ :*

$$(1) f(x) = f(y) + 1 \Leftrightarrow y \preceq x;$$

$$(2) f(x) = 0 \Leftrightarrow x = \perp.$$

The potential benefit of this metric is for direct the allocation. We can update the typical resource allocation algorithms, such as the First Fit and Best Fit by adding this metric to direct the search or use it as a eigenvalue to design supervised learning methods.

**Complexity analysis:** This ranking function is easy to compute on any modular lattice in polynomial time complexity, or say within  $O(n)$  in which  $n$  is the amount of RAr. Based on a greedy approach, we can first define the ranking of a empty RAr to be zero as it must be a bottom in the topology and label the following nodes through iterative search in which we labeled by adding one to each child node from its parent node iteratively. This is by the case that within the result modular lattice, one node will connect with at most two nodes in one direction. Recall the complexity of implementing the representation process is also within polynomial time, then we can say the complexity of the whole process to compute ranking over arbitrary resource allocation topology is in polynomial time.

Note that it is in NP-hard to compute such a ranking function even there is

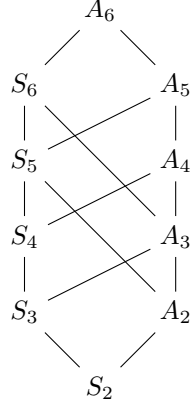


Figure 7: MPC-X RAr topology

in a general poset topology which is as same as solving *topological sorting* problem. We cannot compute all the possible different linear extensions of a typical poset in polynomial time, so cannot label them as well. Then from a complexity point of view, another achievement our approach for ranking the general RAr topology reduce a NP problem into P-time scenario.

## 7 Evaluation

We illustrate some basic experimental evaluations in this section. Say that we only implement some basic prototype programs at this step of the project. The whole program implementation by the novel model illustrated in this report refer to the future work of this two-term project and we describe a plan in the final subsection.

### 7.1 Evaluation Architecture

Two prototype programs are implemented in this step. The first one is a one-dimensional bin packing program by basic linear programming approach. It contains a tiny input size composed by 3 balls and 2 bins. We also implemented a extension vision of this model to handle two-dimensional bin packing problem which is trivially extended from the first one by parallel adding the second dimension. It contains a little larger input size compared with the first one which we used as the basic input instance for multi-dimensional pin packing program.

Our new model is implemented from the third prototype program which is a one-dimensional bin packing program with the same tiny input size of program one, but the allocation is represented by a lattice which using the domain of bins to represent the balls and the values in the domain to represent the size of each bin.

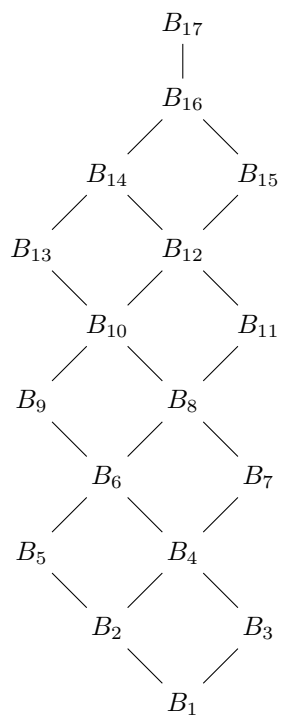


Figure 8: MPC-X RAr lattice-topology

Statistics Results by Novel Model							
No.	Build- ing time	Initial propa- gation	Reso- lution	Nodes	Back- tracks	Max depth	Cons- traints
1	0.0053s	0.010s	0.016s	22	20	21	4
2	0.0052s	0.009s	0.016s	22	20	21	4
3	0.0051s	0.010s	0.017s	22	20	21	4
4	0.0047s	0.010s	0.017s	22	20	21	4
5	0.0065s	0.010s	0.018s	22	20	21	4
6	0.0054s	0.009s	0.016s	22	20	21	4

We apply the same well-known greedy search algorithm as a benchmark of our following experiment.

All of the programs are implemented in Choco3 solver [16] by Java in Eclipse IDE (Vision: Luna Service Release 1a (4.4.1)). The experimental environment is in a MacBook Pro laptop with 2.8GHz Intel Core i7 CPU; 8 GB 1600 MHz DDR3 Memory; and OS X 10.9.2 Operation System.

## 7.2 Evaluation Results

We first compare the programs through trivial linear-programming bin packing model and our novel model. We show our solutions in the two tables which indicate the statistics results in trivial model and our novel model, respectively. When both of them find the same solution, the model’s tuning-time related metrics: building time, Initial propagation and resolution used for program by novel model are all less than the program by trivial model. In addition, the number of constraints applied in our novel model are less than in trivial model.

The interesting results we find is that our model obtains the efficient solutions although it actually suffers more intern nodes and backtrack and reach larger max depth during the searching process through the same search algorithm. It is corresponding to the fact that the model by lattice when always search all the internal subset during the path even it cannot be the solution.

## 7.3 Future Evaluation Plan

We intend to implement our model and the ranking function mentioned in section 6 for the MPC-X device in the next term. Instead of just evaluating the relation-oriented lattice model with the traditional linear programming model through the same trivial searching algorithm as we did above, we will implement the corresponding ranking function for the MPC-X topology and add it as a metric in different searching algorithms to compare the performance with the traditional solver. In detail, we plan to add the ranking function as a metric to

Statistics Results by Trivial Model							
No.	Build- ing time	Initial propa- gation	Reso- lution	Nodes	Back- tracks	Max depth	Cons- traints
1	0.0076s	0.013s	0.020s	4	2	3	5
2	0.0076s	0.014s	0.019s	4	2	3	5
3	0.0083s	0.010s	0.015s	4	2	3	5
4	0.0077s	0.012s	0.017s	4	2	3	5
5	0.0077s	0.011s	0.016s	4	2	3	5
6	0.0078s	0.011s	0.016s	4	2	3	5

update the Best Fit and First Fit searching algorithms to direct the allocation and compare with the non-updated cases. Real-world experimental dataset of RArS to MPC-X device will be used as the input benchmark during the evaluation.

## 8 Conclusions and Future work

This project provides a novel approach to model the resource allocation problem on heterogenous distributed systems by focusing on the relations between resource allocation requests (RArS) and servers. Modelling the resource allocation process as bin-packing problem has been common used in the literature and even viewing this underlining problem from the set theory point of view has been discussed in the pervious work. The novel feature of our work is originally applying Birkhoff’s representation theorem from order and lattice theory to represent the topology of RAr as a modular structure which is much more efficient compared with just modelling them as sets and building up constraints through the values in domain. Instead of using linear-programming approach to solve the constrained optimisation problem, our model is constructed based on representation, or say mapping between the structures of RArS and cluster of servers. In this report, we have detailed illustrates the RAr models from the basic homogenous case to the general heterogeneous and practical case. We invoke ranking function on the lattice structure to give metric on the topology and compute the difference between RArS to direct the allocation.

The main achievement of this project is that we successfully model general heterogenous distributed systems, such as the MPC-X device to a modular lattice instead of a general unconstrained structure. Then we can give easy but efficient metric on RArS to such systems. Our novel model shares a new light on handling resource allocation problem and the underlining multi-dimensional bin-packing problem and specially suited for the case when there are multi-servers RArS and constraints between heterogenous servers.

We have evaluated the efficiency of our approach in the basic case by solv-

ing a trivial bin-packing problem while the implementation of a complicate programme to solve a real resource allocation problem through our model and compare the performance with traditional methods remains to be the future work. Other interesting extensions include applying our model on different practical distributed systems and HPC devices. As our model give a new metric on general RARs, it will also be possible to apply it on an Information Retrieval programme by defining efficient information order and distance.

The main limitation of our model currently is that it doesn't solve the work-time issue: Our model is based on the bin-packing underlining problem instead of a job-scheduling scenario. Recent work such as [24] model the resource allocation as job-scheduling to handle the work-time issue while our model focuses on the static scenario. It is possible to extend our model to solve job-scheduling problem as well but updating the relations and topology between different RARs may thoroughly reduce the efficiency of this model.

## 9 Appendices

### 9.1 VP\_CPSUM Algorithms

We give high level descriptions of VP\_CPSUM algorithm in this subsection as it is the most efficient methods for the basic homogenous model. We first illustrate the VP\_PPSUM.<sup>13</sup> below as it provides the basic structure for designing VP\_CPSUM algorithm.

#### VP\_PPSUM Algorithm

- 1.place each of the  $N$  vectors in on of  $d! / (d-\omega)!$  Lists, where  $\omega \in [1, \dots, d]$   
 $\Rightarrow$  Each list contains the vectors with a common permutation of their largest  $\omega$  dimensions.
- 2.Vectors in each list are then sorted according by decreasing sum of the coordinates (SUM).
- 3.Filling bins with vectors, each time attempting to reduce the resource load imbalance in a bin.  
 $\Rightarrow$  When  $\omega=2$ , first look in list  $(i,j)$  for the first vector that can fit in the bin, hope to reduce the resource imbalance.  
 $\Rightarrow$  If not, relaxes the ordering of the components and searches in other lists.  
 $\Rightarrow$  If no vector can fit in the current bin, then a new bin is added and the process is repeated until all vectors are placed in bins.

#### VP\_CPSUM Algorithm

A relaxation of the VP\_PPSUM algorithm in that it does not enforce any ordering between the  $\omega$  coordinates of vectors and thus need only  $d! / \omega! (d-\omega)!$

---

<sup>13</sup>VP\_PPSUM: Vector Packing algorithm using permutation packing and sorting ordered by decreasing sum of corrdinates

Lists<sup>14</sup>.

### Bin-centric View

Rina Panigrahy, Kunal Talwar, Lincoln Uyeda, and Udi WiederFrom in [15] illustrate the same algorithm<sup>15</sup> from a bin-centric view: Open only one bin at any time, place items into this bin from the largest suitable one. Close the bin when no item can be put in. Then there are two heuristics by involving random choosing. Grasp[k]: pick a random one from the best k instead of the best one; Bubblesearch: the kth best is chosen with a propobability proportional to  $(1 - p)^k$ , for a suitable p.

### Algorithm Analysis

For fixed d and  $\omega$ , both algorithm have complexity  $O(n \log n)$ . In addition, as analysis in [15], for the case when some demands in a dimension always dominate the demands in the other dimension, VP\_CPSUM is more robust. The dimensions which are not scare can be assigned smaller coefficients and have a smaller impact on the ordering.

## 9.2 Heuristics in Wrasse Solver

This subsection illustrated some heuristics used during the searching processing in the Wrasse solver [17]. Those materials are interested for designing efficient resource allocation algorithm but not closed to the modelling process so we keep them in the Appendices.

Wrasse invokes the user defined *dynamic* utilization function with the partial assignment. In briefly: first checks that capacity constraints are met. After that assign friend balls, check for group-consistency, if violated then left them as unallocated. When no more balls can fit, the solver mores on to the next bin.

As for the generic search algorithm it used: When picking balls and bins, different choices optimize for difference outcomes. To balance utilization across bins, it uses power of two sized bin-groups, exploit real-world spatial coherence in bins while being domain agnostic. Wrasse solver tries several size in parallel instead of requiring the user to specify in this process and searched the space by exploration: each thread-group independently exploring its own partial solutions as deep as possible. There were four reasons for choosing explore rather than exploit in the wrasse solver: 1. Define 'good' partial solution in a domain-agnostic manner is problematic; 2. We can retain early lucky decisions; 3. After several steps, we end up with a larger diversity of assignments; 4. Do not need addition heuristics.

---

<sup>14</sup>Giving  $\omega=2$  leads to good empirical result has been shown in [10]

<sup>15</sup>They named their algorithm as FFDAvgSum which means First Fit by Decreasing Average Sum and is the same as VP\_CPSUM in [19]



Analysis by evaluation: For the VM placement problem, Wrasse is able to match the performance of the carefully tuned, production level placement heuristics. For the network virtualization: Both heuristics find the smallest sub-tree in the network that can for a Virtual Cluster. Wrasse’s parallel search for allocations provides high-quality solutions within a few seconds.

### 9.3 User Guide

As illustrated in above section, we use Java-Choco [16] solver and Eclipse IDE to implement the experimental programme. We provide a concise guideline in this section about how to install and use it.

The Java 1.8 can be downloaded from its official website by Oracle through the following URL:

*[https : //www.java.com/en/download/](https://www.java.com/en/download/)*

The Eclipse IDE can be downloaded from the URL:

*[https : //eclipse.org/downloads/](https://eclipse.org/downloads/)*

The Choco solver should be added in the project’s referenced libraries and the jar file we used is *choco-solver-3.3.3-with-dependencies.jar* which can be free downloaded from the URL:

*[http : //choco – solver.org/Download?q = releases](http://choco-solver.org/Download?q=releases)*

The source codes of our prototype programs are open and can be downloaded from the project account in Imperial College’s Gitlab:

*[https : //gitlab.doc.ic.ac.uk/ty215/MResProjectIC.git](https://gitlab.doc.ic.ac.uk/ty215/MResProjectIC.git)*

## References

- [1] Microsoft Assessment and Planning Toolkit (MAP). <http://www.microsoft.com/map/>.
- [2] Hitesh Ballani, Paolo Costa, Thomas Karagiannis, and Ant Rowstron. Towards predictable datacenter networks. In *ACM SIGCOMM Computer Communication Review*, volume 41, pages 242–253. ACM, 2011.
- [3] Nikhil Bansal, Alberto Caprara, and Maxim Sviridenko. Improved approximation algorithms for multidimensional bin packing problems. In *Foundations of Computer Science, 2006. FOCS'06. 47th Annual IEEE Symposium on*, pages 697–708. IEEE, 2006.
- [4] Daniel Bleichenbacher and Jiirg Schmid. Computing the canonical representation of a finite lattice. In *Semantics of programming languages and model theory*, volume 5, page 269. CRC Press, 1993.
- [5] Alberto Caprara and Paolo Toth. Lower bounds and algorithms for the 2-dimensional vector packing problem. *Discrete Applied Mathematics*, 111(3):231–262, 2001.
- [6] FP7 HARNESS consortium. The harness platform: A hardware- and network-enhanced software system for cloud computing. Technical report.
- [7] Brian A Davey and Hilary A Priestley. *Introduction to lattices and order*. Cambridge university press, 2002.
- [8] George Grätzer. *Lattice theory: foundation*. Springer Science & Business Media, 2011.
- [9] Chuanxiong Guo, Guohan Lu, Helen J Wang, Shuang Yang, Chao Kong, Peng Sun, Wenfei Wu, and Yongguang Zhang. Secondnet: a data center network virtualization architecture with bandwidth guarantees. In *Proceedings of the 6th International COnference*, page 15. ACM, 2010.
- [10] William Leinberger, George Karypis, and Vipin Kumar. Multi-capacity bin packing algorithms with applications to job scheduling under multiple constraints. In *Parallel Processing, 1999. Proceedings. 1999 International Conference on*, pages 404–412. IEEE, 1999.
- [11] Siva Theja Maguluri, R Srikant, and Lei Ying. Heavy traffic optimal resource allocation algorithms for cloud computing clusters. *Performance Evaluation*, 81:20–39, 2014.
- [12] Heikki Mannila and Christopher Meek. Global partial orders from sequential data. In *Proceedings of the sixth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 161–168. ACM, 2000.

- [13] Hien Nguyen Van, Frederic Dang Tran, and Jean-Marc Menaud. Autonomic virtual resource management for service hosting platforms. In *Proceedings of the 2009 ICSE Workshop on Software Engineering Challenges of Cloud Computing*, pages 1–8. IEEE Computer Society, 2009.
- [14] GALib: A C++ Library of Genetic Algorithm Components. <http://lancet.mit.edu/ga>, 2010.
- [15] Rina Panigrahy, Kunal Talwar, Lincoln Uyeda, and Udi Wieder. Heuristics for vector bin packing. *research.microsoft.com*, 2011.
- [16] Charles Prud’homme, Jean-Guillaume Fages, and Xavier Lorca. *Choco3 Documentation*. TASC, INRIA Rennes, LINA CNRS UMR 6241, COSLING S.A.S., 2014.
- [17] Anshul Rai, Ranjita Bhagwan, and Saikat Guha. Generalized resource allocation for the cloud. In *Proceedings of the Third ACM Symposium on Cloud Computing*, page 15. ACM, 2012.
- [18] Paul Shaw. A constraint for bin packing. In *Principles and Practice of Constraint Programming–CP 2004*, pages 648–662. Springer, 2004.
- [19] Mark Stillwell, David Schanzenbach, Frédéric Vivien, and Henri Casanova. Resource allocation algorithms for virtualized service hosting platforms. *Journal of Parallel and Distributed Computing*, 70(9):962–974, 2010.
- [20] Mark Stillwell, Frédéric Vivien, and Henri Casanova. Dynamic fractional resource scheduling versus batch scheduling. *IEEE Transactions on Parallel and Distributed Systems*, 23(3):521–529, 2012.
- [21] Mark Stillwell, Frederic Vivien, and Henri Casanova. Virtual machine resource allocation for service hosting on heterogeneous distributed platforms. In *Parallel & Distributed Processing Symposium (IPDPS), 2012 IEEE 26th International*, pages 786–797. IEEE, 2012.
- [22] Maxeler Technologies. *MPC-X*, <https://www.maxeler.com/products/mpc-xseries>. Maxeler Tech., 2015.
- [23] Bhuvan Urgaonkar, Prashant Shenoy, Abhishek Chandra, Pawan Goyal, and Timothy Wood. Agile dynamic provisioning of multi-tier internet applications. *ACM Transactions on Autonomous and Adaptive Systems (TAAS)*, 3(1):1, 2008.
- [24] Shahin Vakilinia, Mustafa Mehmet Ali, and Dongyu Qiu. Modeling of the resource allocation in cloud computing centers. *Computer Networks*, 91:453–470, 2015.
- [25] Hien Nguyen Van, Frederic Dang Tran, and Jean-Marc Menaud. Sla-aware virtual resource management for cloud infrastructures. In *Computer and Information Technology, 2009. CIT’09. Ninth IEEE International Conference on*, volume 1, pages 357–362. IEEE, 2009.