

Domain sizing in Optical Traffic Grooming based Data Center Networks

Ganesh C. Sankaran^{1,2,3} and Krishna M. Sivalingam^{1,3}

¹Department of Computer Science and Engineering, Indian Institute of Technology Madras, Chennai, INDIA

²HCL Technologies Ltd, Chennai, INDIA

³India-UK Advanced Technology Centre of Excellence in Next Generation Networks, Systems and Services (IU-ATC)

Email: gsankara@hcl.com, skrishnam@iitm.ac.in, krishna.sivalingam@gmail.com

Abstract—Optically groomed data center network (OGDCN) is a hybrid optical data center network. This network uses collision domains and wavelength division multiplexing to interconnect all compute storage nodes (CSN). The collision domain size can be adjusted by reconfiguring optical circuit switch (OCS). Two OCS configurations Extender and Isolator are considered for expanding and shrinking the domain size respectively. A Mixed integer linear programming (MILP) formulation is presented for placement of these configurations across the network given the traffic demands. Two variants single domain mapping (SDM) and multiple domain mapping (MDM) of the formulation are presented. These variants are evaluated. It is observed that the MDM variant is better suited for network design and the simpler SDM variant can be used for computing reconfigurations. Also with 5000 flows, it is shown that OCS reconfiguration is not required with a change in traffic distribution without any change in load.

I. INTRODUCTION

Data center networks aggregate compute and storage requirements for an enterprise. Cloud service providers host compute and storage for multiple customers. These networks consume significant amount of power. To address power consumption, many hybrid network architectures [1], [2], [3] were proposed. However, throughput of these architectures were primarily limited by the throughput of Top of the rack (ToR) switches. To address this optically groomed data center network architecture (OGDCN) was proposed in [4].

Every switch in this OGDCN architecture has an optical switch fabric that uses multiplexer, demultiplexer and a set of optical circuit switches. With this fabric, optical transparency is achieved across all network paths. This network is completely agnostic of the data rate or encoding used for communication and can support different data rate and encoding variants simultaneously.

Within the switch fabric, a demultiplexer groups the incoming light signals into wavelength sets. An optical circuit switch is connected to every wavelength set. The input-output connectivity of every wavelength set is determined by the corresponding optical circuit switch configuration. This circuit switch connectivity is critical and determines the wavelength reuse factor across the network. For instance, when wavelength reuse factor of one is desired, all switch fabrics across the network are configured such that all wavelengths form a single large collision domain. The corresponding network throughput supported by this configuration is the minimum possible with this architecture. When higher wavelength reuse factor is

desired, a wavelength set is reserved for communication within the span of the switch, thus creating many smaller collision domains across the network. This increases the overall network throughput. In the current work, the set of sources and destinations interconnected using a wavelength constitute a collision domain.

Roughly, this scenario is more or less similar to macro, micro, pico and femto cell deployment. When using one macro cell and without any smaller cells, the throughput is minimum. With more pico and femto cells, the same spectrum can be reused across the micro cell. This increases the network throughput.

The problem of determining the configuration for a given traffic demand is called the Domain sizing problem. This problem is similar to a bin packing problem with multiple bins. The traffic flows are mapped to articles and domains to bins that are used to pack these articles. An integer linear programming (ILP) formulation is presented to solve this problem. Two variants of this formulation are considered. The first variant single domain mapping (SDM) maps a traffic flow to a collision domain and computes only the wavelength assignment for a flow. The second variant multiple domain mapping (MDM) considers a set of candidate collision domains for every flow. It computes both domain selection and wavelength assignment for the input flows.

These variants are evaluated to determine the trade-off between optimality and resource complexity. An optimal configuration thus determined is deployed on the network. After deployment, the actual input traffic demand can vary. It need not remain same as the input traffic demand. So, the ability to tolerate variations in traffic distribution is evaluated for a smaller set of configurations.

From the evaluation, it is observed that both variants achieve throughput but the number of domains used by SDM variant is more than MDM variant for the same demand set. Thus, SDM variant is more suited for quick network reconfiguration owing to its simplicity where as MDM variant can be used during elaborate design evaluations. A subset of configurations were evaluated for their tolerance to variation in traffic distribution. It was observed that these configurations were found to support all demand sets. Thus reconfiguration frequency can be reduced unless there is variation in network load.

Rest of the work is organized as follows: In Sec. II, necessary background on optically groomed data center networks to understand the current work is presented. Domain sizing

is limited. Thus collision domains must distribute transmitters and receivers such that the flow bandwidth can be supported by the domain it is assigned to. Network throughput variations can be supported by network reconfigurations. When more network throughput is needed, circuit switches across the network can be reconfigured to create more collision domains.

In this section, the background required to understand the rest of the work was discussed. Having had an intuitive understanding of how domain configuration can impact the network throughput, per transmitter and per receiver bandwidths, the next section presents the formal definition of the Domain sizing problem.

III. PROBLEM STATEMENT

In this section, mixed integer linear programming (MILP) formulation of the Domain sizing problem is presented.

The objective is to compute flow selection (a_{ijk}) that maximizes network throughput for the given flow demands and their candidate domain mapping. Apart from these inputs, link capacity R and number of wavelengths m are also specified.

To formulate the problem as a mixed ILP problem, following notations are used.

- 1) Let $F = \{f_1, f_2, \dots, f_n\}$ be the set of flow demands. Here let f_i denote the flow rate.
- 2) Let $W = \{\lambda_1, \lambda_2, \dots, \lambda_w\}$ be the set of wavelengths.
- 3) Let $G = \{g_1, g_2, \dots, g_m\}$ be the set of domains in the network.
- 4) Let a_{ijk} indicate whether the flow f_i is assigned to the domain g_j on wavelength λ_k .
- 5) Let b_{jk} indicate whether the domain g_j is selected on wavelength λ_k for flow assignments.
- 6) Let d_{jk} indicate the cumulative rate of the domain g_j on wavelength λ_k . Thus, $d_{jk} = \sum_i f_i a_{ijk}$, $\forall g_j \in G$ and $\lambda_k \in W$.
- 7) Let the input c_{ij} indicate whether a flow f_i can be assigned to the domain g_j .
- 8) Let R be the capacity of a domain. It is assumed to be constant across the network.

For convenience, suffix i , j and k are used to refer to a specific flow, domain and wavelength within their respective sets.

F , G , W and c_{ij} are the inputs. Two *Boolean* selection variables capture the output. They are a_{ijk} and b_{jk} .

$$\begin{aligned}
 & \text{maximize} && \sum_j \sum_k d_{jk} \\
 & \text{subject to} && a_{ijk} \leq c_{ij}, \quad \forall f_i \in F, g_j \in G, \lambda_k \in W \quad (1) \\
 & && \sum_i \sum_j \sum_k a_{ijk} \leq 1 \quad (2) \\
 & && \sum_{g_j \in P(t)} b_{jk} \leq 1, \quad \forall \lambda_k \in W, \forall t \in T \quad (3) \\
 & && \sum_j \sum_k d_{jk} \leq b_{jk} R \quad (4)
 \end{aligned}$$

Here T is the set of all ToR switches and $P(t)$ is the set of domains along the network path from the ToR switch t to any core switch on the network.

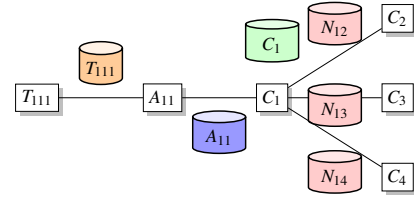


Fig. 4: Domain selection constraint set $P(t)$ for a specific ToR switch T_{111}

Eqs. 1–4 present flow to domain mapping, flow selection, domain selection and domain capacity constraints respectively. The flow to domain mapping constraint ensures that a flows selection within candidate domains. The flow selection constraint ensures that a flow is mapped on to only one domain across the network. The domain selection constraint ensures that one wavelength from any compute/storage node is mapped to only one domain along the network path from the ToR to the network core as discussed later. The domain capacity constraint ensures that sum of flow rates for every selected domain is within the capacity. Except the flow and domain selection constraints, the other constraints are from the multiple bin packing problem.

The flow selection constraint ensures that only one domain and one wavelength is assigned to an input flow. Similarly, domain selection constraint enables selection of valid domains on all supported wavelengths. One domain selection constraint considers the path from a ToR switch (t) to all the core switches in the network. This path stops just ahead other core switches and thus it does not include other core switches except for the one that includes the ToR switch t within its span. The set of domains appearing along these paths strands is represented by the domain set $P(t)$. A mutual exclusion condition is enforced on these domains. Thus only one domain can be selected on a wavelength within this domain set. For instance, $P(T_{111}) = \{T_{111}, A_{11}, C_1, N_{12}, N_{13}, N_{14}\}$ is shown in Fig. 4.

A close observation reveals that domain selection constraint considers only the one half of the network path starting from a ToR switch to the core switches. This mutual exclusion of domains nicely extends itself to other half network path. To illustrate this, let us consider the network path from T_{111} to any core switch say C_2 . This is shown in Fig. 5. When N_{12} domain is selected on specific wavelength say λ_k the domain selection constraint can be satisfied only when none of the other domains (C_1 or A_{11} or T_{111}) in the set are selected. In this case, the other half of the network path is not covered by this constraint. However when collectively all constraints involving N_{14} are enforced, it ensures that no other domains are selected within the span of core switches C_1 and C_2 on wavelength λ_k .

Similarly, domain span can be restricted to a specific tier of the network. For instance, a domain can include the network span within a core, aggregate or a ToR switch. When selecting all domains within the span of a ToR switch, the number of domains increases and so does the network throughput. However, domain size decreases in this case. When selecting

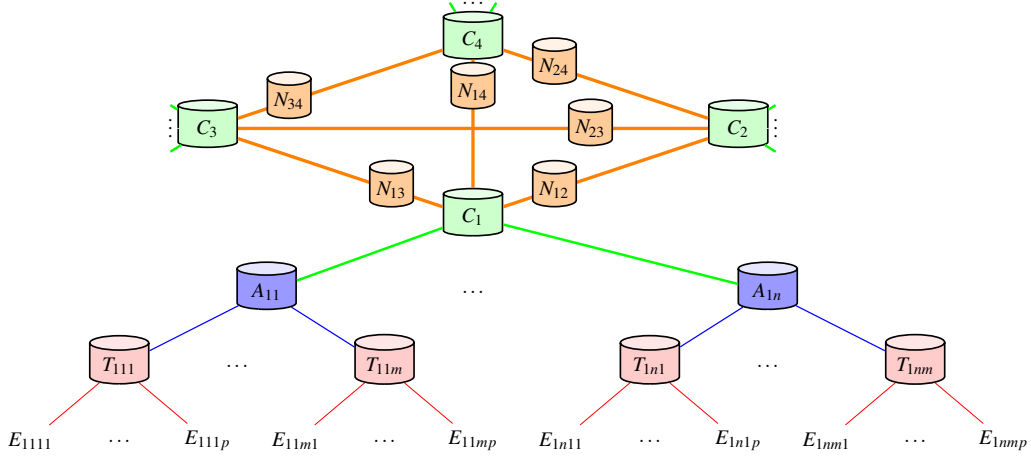


Fig. 5: Domains represented by corresponding bins

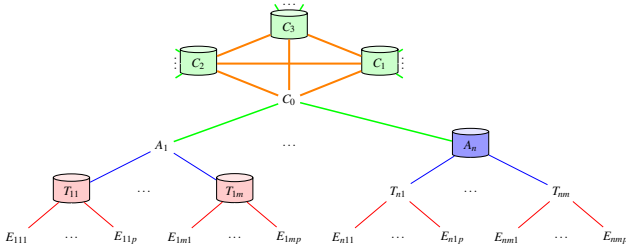


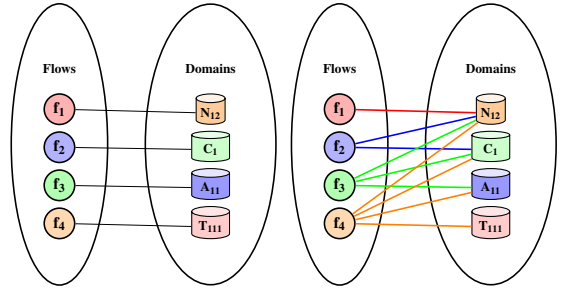
Fig. 6: A realistic domain selection across the network

all domains within the span of core switches, the number of domains decreases but the domain size increases. This results in reduction in network throughput. An example domain selection across network tiers is shown in Fig. 6.

This problem is similar to multiple bin packing [5] problem where multiple bins are packed with articles of different weights such that the overall weight of the bins after packing is the maximum. The multiple bin packing problem can be mapped onto the current problem by mapping bins to domains, articles to flows, article weights to flow rates.

While multiple bin packing problem and Domain sizing problem are similar so far, there are key differences. First, a flow cannot be mapped onto any arbitrary domain, so flow selection constraints restrict the flow to domain mapping. For instance, a flow within the span of the core switch C_1 cannot be mapped onto a domain within the span of another core switch. Second, all domains cannot be selected simultaneously. One wavelength from a compute/storage node can be mapped onto only one domain. So domain selection constraints exist that aid selection of a valid set of domains across the network. Thus the Domain sizing problem is a more constrained version of the multiple bin packing problem.

A domain corresponds to a network path and is denoted by the network equipment in the top-most tier along the network path. For example, the domain that encompasses the network within the span of the core switch C_1 is denoted as C_1 . One exception to this domain notation is the network-wide domains that contain two core switches. These domains are denoted by N_{xy} where x, y are suffix of the corresponding core



(a) One to one mapping SDM (b) One to many mapping MDM

Fig. 7: Flow to Domain Mapping: One to one (SDM) and one to many (MDM) mapping variants. Flow f_1 uses a network path that contains two core switches C_1 and C_2 . Flows f_2 , f_3 and f_4 use a network path within the span of a core, aggregate and ToR switch respectively.

switches. For example, N_{14} denotes the network-wide domain that encompasses the sub-networks within the span of core switches C_1 and C_4 . It also includes the link between these core switches. The significance of network-wide domains is discussed later in the context of domain selection constraints.

The variants to be evaluated with this formulation are introduced.

SDM Variant: This variant considers flow to domain mapping as one-to-one. Thus the formulation is reduced to finding the wavelength to be assigned for this domain as shown in Fig. 7(a). Thus any given flow the value of c_{ij} is set only for one domain across all domains. This is set to the domain that is nearest to the source. For instance, flow f_4 can be supported by more than one domain. One among the domains denoted by T_{111} or A_{11} or C_1 or N_{12} or N_{13} or N_{14} can be chosen. In this variant, the domain nearest to the source is T_{111} so this is chosen. Considering another flow f_1 , it requires a network path that passes through core switches C_1 and C_2 . This path requirement is satisfied by only one domain N_{12} . Thus f_1 is mapped onto domain N_{12} . Similarly, flows f_2 and f_3 are mapped onto A_{11} and C_1 respectively.

MDM Variant: This variant considers flow to candidate

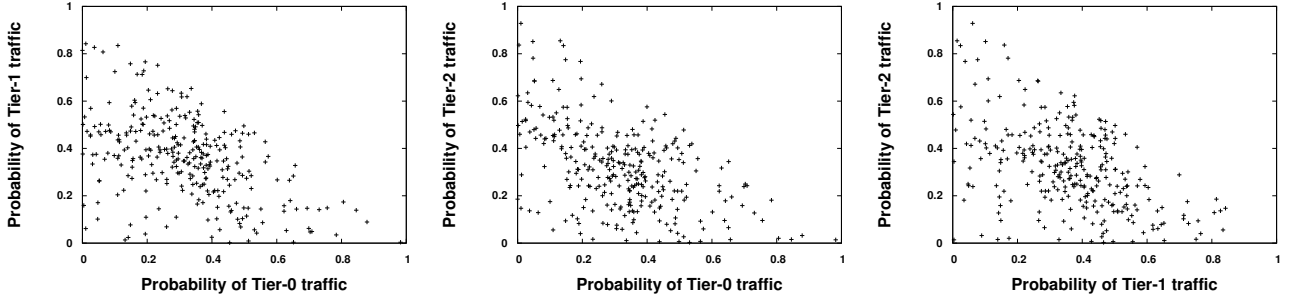


Fig. 8: Probability distribution of flows

domain mapping as one-to-many. Thus the formulation is required to find one of these domains and a wavelength for the selected domain to meet the objective. Thus with respect to flow assignment this formulation has two degrees of freedom (one domain and another wavelength). This mapping is shown in Fig. 7(b). In this variant, all domain choices are provided to the formulation and the domain that meets the objective can be chosen.

The primary goal of defining these variants is to analyze the solution quality of these variants to see whether providing additional degree of freedom in terms of choice of domains results in better solution quality.

In this section, the problem formulation and the proposed variants were discussed. The problem was shown to be a constrained version of the multiple bin packing problem. Two variants SDM and MDM were defined. These variants will be evaluated in the next section.

A. Configuration example

When a domain say C_1 is selected on a wavelength λ_k , the corresponding network configuration for the wavelength is as follows. The core switch multiplexer and demultiplexer will redirect the wavelength λ_k through Isolator OCS configuration. All other switches within the span of the core switch C_1 must pass wavelength λ_k through the Extender OCS switch configuration. This will result in a domain that includes all network paths within the span of the core switch C_1 and excludes other core switches in the network. This domain can interconnect all compute/storage nodes within the span of the core switch C_1 .

IV. EVALUATION

In this section, the SDM and MDM variants that were defined earlier are evaluated.

To evaluate the variants, following setup and parameters were used. A server with 32 cores each clocking 2 GHz with 24 GB RAM was used for solving the MILP formulations in parallel. A set of 300 traffic demands were generated. Every demand is made up of 20,000 flows. The flow rate for every flow is picked from an uniform distribution between 0 and 100 Mbps. A network-wide scaling factor is used to set the link or collision domain capacity to either 1 or 10 Gbps. The number of wavelengths used is 64. A two-tier hierarchical topology is considered. The total number of core and ToR

switches across the network is four and 64 respectively. The number of compute/storage nodes in this DCN is 1024.

All domains that can route the flow to its destination are considered as candidate domains for a given flow with the MDM variant. Of these domains, the one that is nearest to the source is its candidate domain with the SDM variant.

The total number of flows are partitioned into three: (i) flows across core switches (Tier-0 flows), (ii) flows within the span of a core switch (Tier-1 flows) and (iii) flows within the span of a ToR switch (Tier-2 flows). The actual number of flows in a traffic demand is based on three random fractions. These fractions are used as the corresponding weight to compute the number of flows in each tier. For instance, when the fraction for Tier-0 flows is 0.3, the number of flows in the corresponding tier is 6000 ($20,000 \times 0.3$).

The distribution of number of flows across tiers for the 300 traffic demands is shown in Fig. 8. In figures (a), (b) and (c) respectively present the fraction of flows in Tier-1 with respect to Tier-0, Tier-2 to Tier-0 and Tier-2 to Tier-1. From the figure, it can be noted that Tier-0, Tier-1 and Tier-2 fractions vary from 0 to 0.98, 0.001 to 0.84 and 0.003 to 0.92 respectively. The variance of these probabilities is between 0.000046 and 0.22. This range is wide and is expected when using an uniform distribution.

Gurobi 6.02 [6], a MILP solver was used to solve the formulation. The throughput achieved, the number of flows allotted, the number of domains used and the maximum domain utilization were measured for every traffic demand.

Scenario 1: This scenario explores the differences between SDM and MDM variants in terms of measured values. Initially, the domain or link capacity was set to 1 Gbps. To vary the load, the number of input flows was increased from 5,000 to 20,000 in steps of 5,000. When the number of flows was 15,000 or more, the solution was found to drop many flows. This indicates that the applied load was beyond its operational range. So, the link or domain capacity was increased to 10 Gbps to accommodate 15,000 or more flows.

From the computation, no throughput difference was observed between SDM and MDM variants for all 300 traffic demands. The same trend is observed with a load of 5,000 and as well as with 10,000 flows. So a simpler formulation can be chosen during reconfigurations to reduce the computational resource requirements. In this case, SDM formulation can be chosen for reconfiguration.

However, there is difference in the number of domains

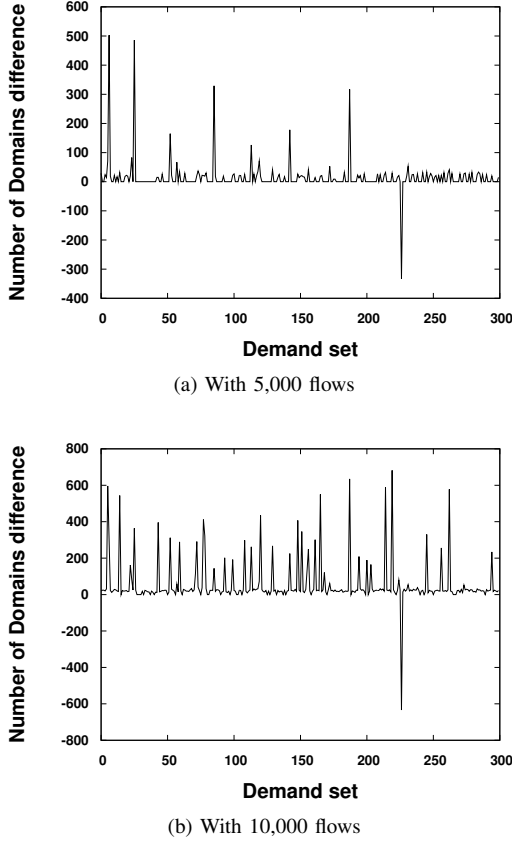


Fig. 9: Difference in number of domains chosen by SDM and MDM variants

Input flows	Demands computed	Full allocation	Least Allotted	Domains used	Max. utilization
5000	300	300	5000	89-1172	99.9%
10000	300	300	10000	88-1176	100%
15000	291	218	14993	89-1138	64%
20000	245	67	19993	119-949	73%

TABLE I: Performance of variants with different flow sizes: Number of demand sets computed within 800 sec. (Demands computed), Number of demands fully allotted (Full Allocation), Minimum number of flows allotted (Least Allotted), Number of Domains used (Domains used) and maximum domain utilization (Max. Utilization)

selected by these variants. This is shown in Fig. 9. From the figure, it can be noted that, MDM variant select less number of domains than SDM variant for all but one traffic demand. This traffic demand has 92% Tier-2 flows out of all input flows. The domain count is significant during the design phase when hardware dimensions are computed. Thus it is observed that MDM variant is suitable for the design phase.

Since the throughput was same for five and ten thousand flows, only the MDM variant was evaluated for more flows.

For 15 thousand flows, 291 demand sets were computable within the 800 second time limit specified for the solver. Of them, all flows in 218 demand sets were completely accommodated. With the rest 73 demand sets, the number of flows was at least 14,993. In other words, at most seven flows

were dropped across all demand sets.

For 20 thousand flows, 245 demand sets were computable within the 800 second time limit specified for the solver. This indicates increase in time complexity with increase in the number of flows. Of all demand sets, 67 of them accommodated all flows. Rest of them were short by at most seven flows. This is same as the number of flows dropped by the formulation with 15 thousand flows.

The ones that did not finish computation within 800 seconds were recomputed without any time limit. This time it took the solver 913 seconds on an average and with a maximum of 1107 seconds. **After solving the MILP, time is required to push the reconfigurations on to the network elements. Thus the minimum reconfiguration interval chosen must be more than the time taken by the solver.**

The number of domains used was between 89 to 1138. The range of domains used was 119-949. This is narrower than the number of domains used for 15 thousand flows.

The maximum domain utilization was 64% and 73% for 15 and 20 thousand flows respectively. With 5 and 10 thousand flows, the corresponding utilization was 99.9% and 100% respectively. This indicates that more flows can be accommodated with the same domain selection across the network by increasing the data rate.

Scenario 2: The collision domains are reconfigured over a long duration, when there is a significant change in the traffic pattern. When traffic pattern changes significantly between successive reconfigurations, the current configuration may not support the change in traffic. This can result in dropping of some flows. To assess the ability of a configuration to support variations in traffic distribution at the same load, this scenario explores the set of traffic demands that can be supported by a random set of domain selections.

Intuitively, when domain selection is optimized for a specific traffic demand is used for another demand, it can drop some flows. At light loads where system is not saturated but still flows can be dropped.

For this, a set of ten randomly selected solutions were chosen. So far the collision domains were computed by the formulation. In this scenario, the collision domains selected by the specific solution is used as an input to the formulation. For this selection, maximum achievable throughput was computed for the 300 traffic demands.

With all 300 traffic demands, these solutions achieved maximum throughput achieved by the optimal configuration. Their capacity was sufficient to support the variations in traffic distribution. It can be observed that a chosen configuration with 5000 flows is capable of supporting variations in traffic distribution. This is despite the fact that the maximum utilization was 99.9% for some domains as observed in the previous scenario. Generalizing the observation a randomly chosen configuration is capable of accommodating traffic distribution variations without increase in the number of flows. Thus, frequent reconfigurations are not required in such a scenario.

Of the randomly selected solutions, the minimum and maximum number of domains observed was 97 and 1163 respectively. The solution that selected 1168 domains used the minimum number of domains to support many traffic demands.

Among the domain selections, the solution that supported minimum number of domains to support a traffic demand was observed. Out of 300 traffic demands, 69, 80, 92 and 97 domains were sufficient for supporting 194 (64%), 1, 1 and 104 (35%) traffic demands respectively. Interestingly, the solution that selected 69 domains was the one that had selected maximum number of domains among the randomly selected domains. The solution that selected 97 domains is the same one that had the minimum number of domains among the randomly selected solutions. Thus there is some correlation between minimum and maximum domain selections and the minimum number of domains. This can be explored further to formulate a suitable non-deterministic heuristic.

So far, the evaluation results were presented. From the first scenario, it was observed that maximum throughput can be achieved by both SDM and MDM variants. Since both formulations compute same throughput, the simpler of the formulations can be used during reconfiguration to reduce computational resources. With the second scenario, the randomly selected solutions were tolerant to wide variations in traffic distribution. The frequency of reconfigurations can be reduced when only traffic distribution changes are observed on the network.

V. CONCLUSIONS

The choice of domains largely influences the maximum achievable throughput. So the domains selected by the solution of the Domain sizing problem can be configured on the network periodically based on changes in traffic demand. To arrive at an optimal domain selection, the Domain sizing problem was defined. This problem was shown to be similar to multiple-bin packing problem. Two variants of the problem, SDM and MDM were defined as mixed integer linear programming problems. The performance of these variants were evaluated. From the evaluation, it was observed that SDM variant is sufficient for computing the domains during reconfiguration. Similarly, it was observed that without reconfiguration, the network is capable of supporting changes in traffic distributions. Thus frequency of reconfigurations can be reduced.

REFERENCES

- [1] G. Wang, D. G. Andersen, M. Kaminsky, K. Papagiannaki, T. Ng, M. Kozuch, and M. Ryan, "c-Through: Part-time optics in data centers," in *Proceedings of ACM SIGCOMM*, pp. 327–338, 2010.
- [2] S. Yoo, Y. Yin, and R. Proietti, "Elastic Optical networking and low-latency high-radix optical switches for Future Cloud Computing," in *IEEE International Conference on Computing, Networking and Communications (ICNC)*, pp. 1097–1101, 2013.
- [3] N. Farrington, G. Porter, S. Radhakrishnan, H. H. Bazzaz, V. Subramanya, Y. Fainman, G. Papen, and A. Vahdat, "Helios: a Hybrid Electrical/Optical switch architecture for modular Data Centers," in *Proceedings of ACM SIGCOMM*, pp. 339–350, 2011.
- [4] G. C. Sankaran and K. M. Sivalingam, "Optical Traffic Grooming based Data Center Networks: Node architecture and Comparison," tech. rep., Indian Institute of Technology Madras, 2015. Under review in *IEEE Journal on Selected Areas in Communications SI: Green Communications and Networking*, <http://www.cse.iitm.ac.in/~skrishnam/OGDCN2014.pdf> Password: jocsch.
- [5] N. Karmarkar and R. M. Karp, "An efficient approximation scheme for the one-dimensional bin-packing problem," in *23rd Annual Symposium on Foundations of Computer Science, SFCS'08*, pp. 312–320, IEEE, 1982.
- [6] "Gurobi optimizer reference manual," 2012. <http://www.gurobi.com>.