

Software-defined optical network for metro-scale geographically distributed data centers

Payman Samadi,* Ke Wen, Junjie Xu, and Keren Bergman

Department of Electrical Engineering, Columbia University, New York, New York 10027, USA

*ps2805@columbia.edu

Abstract: The emergence of cloud computing and big data has rapidly increased the deployment of small and mid-sized data centers. Enterprises and cloud providers require an agile network among these data centers to empower application reliability and flexible scalability. We present a software-defined inter data center network to enable on-demand scale out of data centers on a metro-scale optical network. The architecture consists of a combined space/wavelength switching platform and a Software-Defined Networking (SDN) control plane equipped with a wavelength and routing assignment module. It enables establishing transparent and bandwidth-selective connections from L2/L3 switches, on-demand. The architecture is evaluated in a testbed consisting of 3 data centers, 5–25 km apart. We successfully demonstrated end-to-end bulk data transfer and Virtual Machine (VM) migrations across data centers with less than 100 ms connection setup time and close to full link capacity utilization.

© 2016 Optical Society of America

OCIS codes: (060.4250) Networks; (060.4510) Optical communications; (200.4650) Optical interconnects.

References and links

1. C. Kachris, K. Bergman, I. Tomkos, *Optical Interconnects for Future Data Center Networks*, (Springer, 2013).
2. Bell Labs, “Metro network traffic growth: an architecture impact study,” Strategic White Paper 1–12, (2013).
3. Cisco, “Cisco visual networking index: forecast and methodology, 2014–2019,” White paper 1–14, (2015).
4. P. Samadi, J. Xu, K. Wen, H. Guan, Z. Li, K. Bergman, “Experimental demonstration of converged inter/intra data center network architecture,” in *Proceedings of International Conference on Transport Optical Networks (ICTON)*, (2015), paper We.B3.3.
5. N. McKeown, T. Anderson, H. Balakrishnan, G. Parulkar, L. Peterson, J. Rexford, S. Shenker, J. Turner, “Open-flow: enabling innovation in campus networks,” *SIGCOMM Comp. Commu. Rev.* **38**(2), 69–74, (2008).
6. J. Zhang, Y. Zhao, H. Yang, Y. Ji, H. Li, Y. Lin, G. Li, J. Han, Y. Lee, T. Ma, “First demonstration of enhanced software defined networking (eSDN) over elastic grid (eGrid) optical networks for data center service migration,” in *Optical Fiber Communication Conference*, (Optical Society of America, 2013), paper PDP5B.1.
7. S. Yan, E. Hugues-Salas, V. J. F. Rancoo, Y. Shu, G. M. Saridis, B. Rahimzadeh Rofoee, Y. Yan, A. Peters, S. Jain, T. May-Smith, P. Petropoulos, D. J. Richardson, G. Zervas, D. Simeonidou, “Archon: a function programmable optical interconnect architecture for transparent intra and inter data center SDM/TDM/WDM networking,” *J. Lightwave Technol.* **33**(8), 1586–1595, (2015).
8. P. Samadi, J. Xu, K. Bergman, “Virtual machine migration over optical circuit switching network in a converged inter/intra data center architecture,” in *Optical Fiber Communication Conference*, (Optical Society of America, 2015), paper Th4G.6.
9. P. Samadi, H. Guan, K. Wen, K. Bergman “A software-defined optical gateway for converged inter/intra data center networks,” *IEEE Optical Interconnect Conference (OI)*, (2015), paper MB4.
10. T. S. El-Bawab, *Optical Switching*, (Springer, 2008).
11. Polatis 384×384 Optical Space Switch, <http://www.polatis.com/>.

12. The OpenDayLight Platform, <https://www.opendaylight.org/>.
 13. RYU SDN Platform, <https://osrg.github.io/ryu/>.
 14. H. Zang, J. P. Jue and B. Mukherjee, "A review of routing and wavelength assignment approaches for wavelength-routed optical WDM networks," *Optical Networks Magazine* **1**, 47–60, (2000).
 15. S. Sanfilippo and P. Noordhuis, "Redis", <http://redis.io>.
 16. Iperf, <https://iperf.fr/>.
 17. S. Han, T. J. Seok, N. Quack, B. Yoo, M. C. Wu, "Monolithic 50x50 MEMS silicon photonic switches with microsecond response time," in *Optical Fiber Communication Conference*, (Optical Society of America, 2014), paper M2K.2.
 18. T. A. Strasser, J. L. Wagener, "Wavelength-selective switches for ROADM applications," *IEEE J. Sel. Top. Quantum Electron.* **16**(5), 1150–1157, (2010).
 19. A. Biberman, H. Lira, K. Padmaraju, N. Ophir, J. Chan, M. Lipson, K. Bergman, "Broadband silicon photonic electrooptic switch for photonic interconnection networks," *IEEE Photon. Technol. Lett.* **23**(8), 504–506, (2011).
-

1. Introduction

Big data and analytics have created a vast demand for small to large scale data centers among various industries including IT, finance and health care. These data centers with 100-1M servers require thorough inter and intra data center connectivity for fast data access and replication [1]. Ideally, servers and users in different locations must obtain access to the data upon generation at the source. In addition, it is desirable to distribute applications over several data centers to improve operation reliability and surpass scalability limits. Hence, on-demand, high bandwidth and low-latency connections are necessary for inter data center networking.

Currently, enterprises and cloud providers deploy several small to mid-sized data centers that are generally located in metro-scale distances. These data centers are actively communicating for services such as data replication, Virtual Machine (VM) migration, backup, load balancing or fault/disaster recovery. Studies predict 5-6 \times increase in total metro network traffic by 2017 with 75% terminating within the metro network [2, 3]. This growth and shift that is mainly IP traffic, along with data center connectivity requirements, necessitates dynamic metro-scale optical networks with seamless communication within layers [4].

Software-Defined Networking (SDN) enables remote management of network components in the data plane from the higher layers. In the context of optical networks, physical layer components and sub-systems become software accessible. Cross-layer designs along with south-bound APIs (OpenFlow [5]) are the enablers and lead to dynamic networks with rapid reconfigurability.

Researchers have identified the importance of dynamic inter data center connectivity and have investigated various approaches to leverage SDN in optical networks. In [6], software defined networking over elastic grid (eGrid) optical networks for data center service migration is proposed. Aragon [7] have leveraged optical space switches as a reconfigurable cross-connect and combined WDM/TDM/SDM to address both inter and intra data center connectivity. However, these methods do not address the need of enterprises for flexible scalability and operation distribution of small and mid-sized data centers to a larger data center based on the applications' and service's demands. In previous works, we introduced the concept of converged inter/intra data center architecture and demonstrated the basic architecture [4, 8].

In this work, we present a software-defined metro network to enable on-demand and bandwidth-selective inter data center connections. The hardware architecture consists of a combined space/wavelength switching platform (optical gateway [9]), enabling both optical wavelength and circuit switching. The optical gateway aggregates racks/pods/clusters of the data center and provides transparent cross data center connections directly from/to L2/L3 switches. It can also transmit east-west traffic by the spare capacity. The software architecture consists of two SDN modules that integrate with the data center and metro network control plane. The metro network module includes a wavelength and routing assignment kernel. The data center

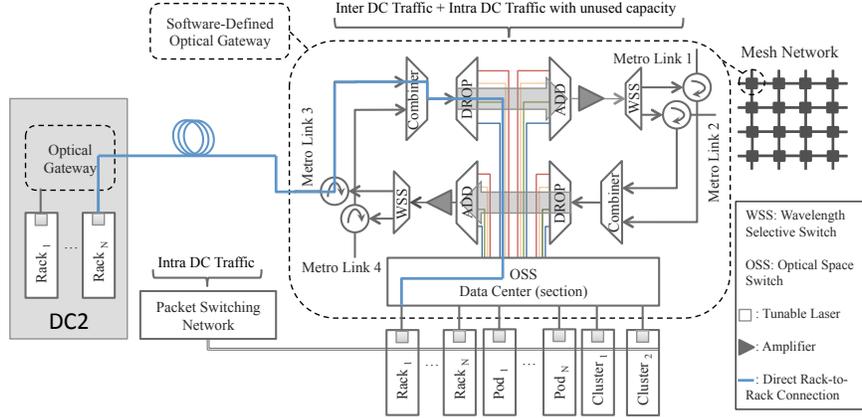


Fig. 1. Hardware architecture of a bidirectional software-defined optical gateway in a mesh network.

module, in a 3-layered architecture, enables cross-layer communication for the physical layer management and also data center to metro network connectivity.

We built a 3-node data center testbed with 5–25 km distance for experimental evaluations. The control plane latency including the optical gateway setup time is measured. For end-to-end evaluations, we demonstrate bulk data transfer and VM migration. In order to verify the scalability of the software architecture, we implemented a simulation platform to assess the wavelength and routing assignment algorithm.

Compared to current methods, our proposed architecture provides following benefits for inter data center networking: i) increasing application reliability by enabling seamless operation distribution over multiple data centers, ii) scaling out data centers in distance to overcome scalability limits and increase total processing power, iii) improving optical link capacity utilization in the metro network by an application driven-approach and finer granularity in managing the wavelength usage, iv) improving networking energy consumption by providing transparent optical links from L2/L3 switches. Majority of enterprises that operate several small to mid-sized data centers can leverage this architecture to manage their own inter data center network through dark fiber and enable bandwidth on-demand based on business requirements.

In the rest of this paper, first we discuss the design of the software-defined optical gateway that is the enabler of the wavelength/space switching platform. In section III, the system architecture is explained. Section IV discusses the testbed implementation. Section V presents the experimental results on the testbed and the numerical results on the scalability of the control plane algorithms. Section VII is a discussion on the future trends of inter data center networking and potential improvements of our architecture and section VIII concludes the paper.

2. Software-defined optical gateway

The optical gateway is the enabler of our proposed space/wavelength switching platform. It performs as a high port count Color-less, Direction-less, Contention-less Reconfigurable Optical Add-Drop Multiplexer (CDC-ROADM) with SDN Interoperability. Figure 1 demonstrates the hardware architecture in a mesh metro network. It consists of an Optical Space Switch (OSS) that is a circuit-based switching substrate [10], a wavelength Selective Switches (WSS) that provides WDM switching, WDM Multiplexers (Mux) and Demultiplexers (Demux) that add/drop WDM channels, and optical combiners and circulators. In the diagram two sets of Mux/Demux and WSS are used to enable bidirectional connections. The control plane of the

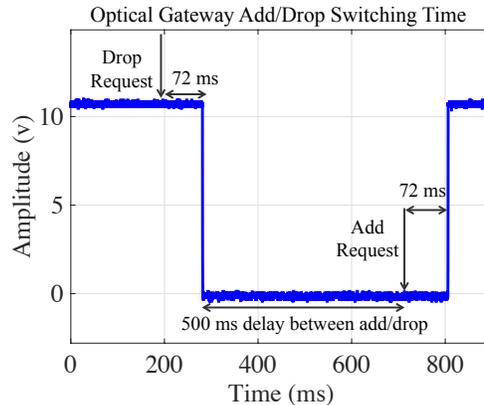


Fig. 2. Switching time of the optical gateway to add and drop a channel, requests are sent sequentially with 500 ms delay in between.

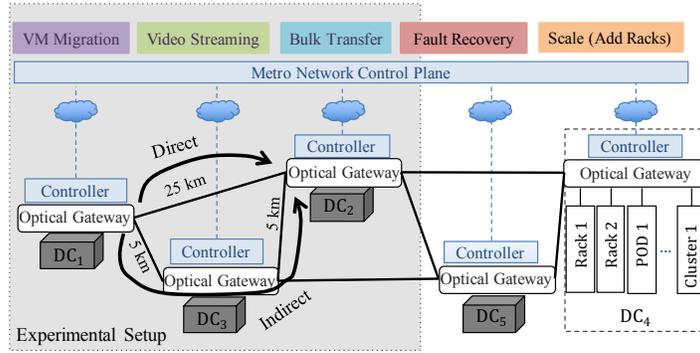
optical gateway is SDN-enabled, i.e. the connectivity configuration is set independent to the data via APIs.

Depending on the scale of the data center and the applications, racks, pods, or clusters of the data center are connected to the optical gateway from the OSS ports. In a sample 4×4 mesh metro network, each node (a data center) is connected to 4 adjacent nodes via Metro Links 1–4. The optical gateway provides the connectivity among the nodes. At each gateway, WDM Mux/Demux selects and combines desired channels. The WSS provides dynamic connectivity of the WDM channels to adjacent nodes. The number of WSS ports is determined by the number of adjacent nodes. In case of a sole connection, the WSS is not required. This architecture scales to support up to 384 racks/pods/clusters using commercially available optical space switches [11].

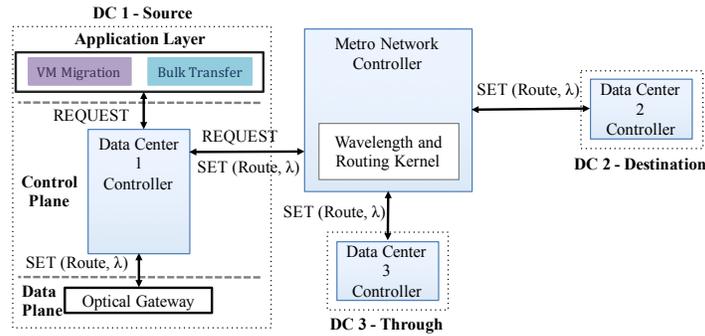
The Electrical to Optical (E/O) conversion is executed by optical transceivers directly on the L2/L3 Top-of-Rack (ToR) switches or aggregation switches of pods or a cluster. We propose leveraging tunable lasers for efficient usage of the wavelength capacity on the optical metro network. Commercial SFP+ transceivers are now supporting 80 km reach on DWDM channels. For longer distances, optical amplifiers are used to boost the output signal to the metro network. In this case, the power balancing of the WDM channels entering the network is performed by the attenuator inside the WSSs.

We employed southbound APIs to enable software control on the optical gateway. We used OpenFlow for the OSS and developed APIs for the WSS and tunable lasers, in-house. All APIs are aggregated in a control plane for add and drop functions. Extending OpenFlow for the WSS and tunable lasers will enable using standard SDN controllers such as OpenDayLight [12] and RYU [13].

We implemented optical gateways using MEMS and piezo-electric optical space switches and Digital Light Processing (DLP) WSS. We measured the switching time of the MEMS-based optical gateway to add and drop a new channel. Request commands to connect rack 1 with wavelength C36 (1548.51 nm) to the Metro Link 4 (Fig. 1), are sent sequentially in a loop with 500 ms delay in between. Figure 2 demonstrates both add and drop switching time that is 72 ms. This is the latency to establish or remove a connection on the physical layer that includes the OSS and WSS switching time but not the power adjustment delays.



(a)



(b)

Fig. 3. (a) Architecture of the software-defined metro-scale inter data center optical network, each data center is equipped with an optical gateway and the metro network control plane manages the connection requests centrally. (b) The control plane workflow to make a connection between DC₁ and DC₂ via DC₃.

3. Architecture

Figure 3(a) demonstrates a metro network consisting of five data centers (DC₁ – DC₅). An optical network provides the connectivity among the data centers in a mesh topology. For distances up to 80 km, commercial SFP+ transceivers provide sufficient power budget. Optical amplifiers (EDFAs) can be used to increase the transmission distance for larger networks. In our architecture, each data center is equipped with an optical gateway that connects to the data center in different tiers such as rack, pod or the aggregation layer (DC₄). The optical gateway provides the north-south connectivity of the data center to the metro network. The spare switching capacity of the optical gateway can be used to transport east-west traffic as well.

The software architecture integrates with the data center and metro network control plane. For the metro network, it consists of i) a controller that sends and receives connection requests from the data centers or the metro network optimizer, and ii) a wavelength and routing kernel that determines the optimal routing and wavelength for the connection requests based on the availability. For the data center, the controller is the connectivity point to the metro network control plane and between the application and data plane layers inside the data center. The connection medium between the data center and the metro network control plane is internet. Secure protocols and firewalls can be used to ensure the security of the control plane. Also, considering the frequency of inter data center connection request, the delay imposed by internet is acceptable.

Figure 3(b) demonstrates the control plane workflow to establish a new connection. Connection requests are made from the application layer for services such as VM migration and bulk data transfer or a network optimizer software. The data center controller receives the requests and forwards it to the metro network. The metro network controller that has the global knowledge of the network including the topology and wavelength availability, assigns a wavelength and a route to the request using the wavelength and routing kernel. Then it sends a set request with the wavelength and route information to the data centers involved in the connection. Each data center controller configures the optical gateway and notifies the metro network controller upon configuration completion. Once the connection is ready, the service starts data transmission.

3.1. Wavelength and routing assignment

The wavelength and routing assignment is implemented using Integer Linear Programming (ILP). The optimization goal here is to minimize the highest indexed wavelength (w_{max}) used to accommodate an incoming connection request. An ILP problem is formulated as follows to jointly optimize routing and wavelength assignment of a connection request. F_{ij}^{sdw} is a binary variable; it has the value 1 when the wavelength w on the link ij is used for the connection from source s to destination d , and 0 otherwise. Eq. (1) is the goal of the ILP, Eq. (2) is the constraint on maximum number of WDM channels. Eq. (3) is the flow conservation constraint at each node and Eqs. (4) and (5) guarantee that the same wavelength is used along the path.

$$\min w_{max} \quad \text{s.t.} \quad (1)$$

$$w_{max} \geq w F_{ij}^{\hat{s}\hat{d}w}, \quad \forall i, j, w \quad (2)$$

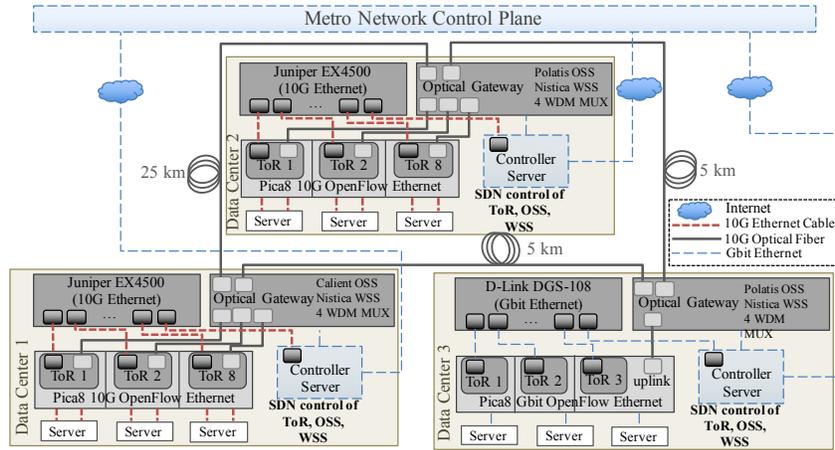
$$\sum_i F_{ij}^{\hat{s}\hat{d}w} - \sum_k F_{jk}^{\hat{s}\hat{d}w} = \begin{cases} -1 & \text{if } j = \hat{s} \\ 1 & \text{if } j = \hat{d} \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

$$F_{ij}^{\hat{s}\hat{d}w} = 0 \text{ or } 1, \quad \forall i, j, w \quad (4)$$

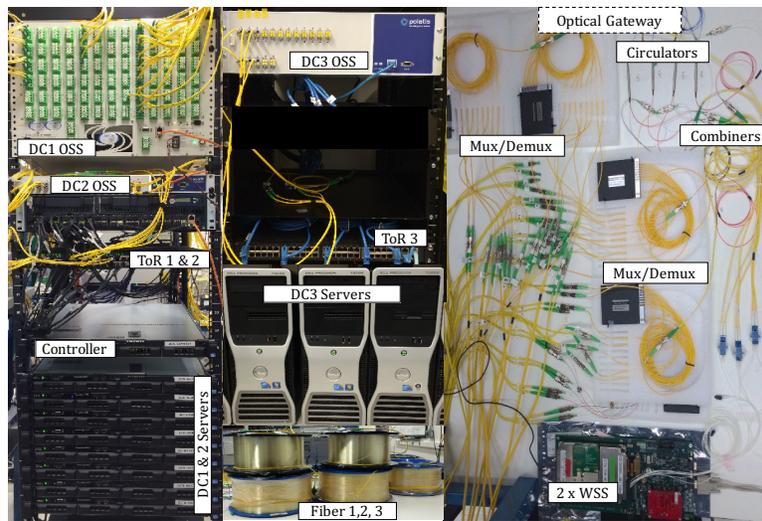
$$\sum_{s,d} F_{ij}^{sdw} \leq 1, \quad \forall i, j, w \quad (5)$$

Finding an optimal solution for the above ILP problem could be of high time complexity and hence prohibitive for real-time applications. Here, we use a *first fit* heuristic to find the path and the wavelength for incoming requests. First fit scheme does not require global information of the network and compared to other methods such as random wavelength, Least-Used (LU) and Most-Used (MU) assignments [14], has lower computation cost since there is no need to search the entire wavelength space for each path request. Furthermore, it performs well in terms of blocking probability and fairness [14].

The heuristic first finds the shortest path between the requested source and destination data centers, then searches for a wavelength that is available on all edges of the path. The search considers a lower-indexed wavelength prior to a higher-indexed one, and the first available wavelength will be selected. If no available wavelength is found, the algorithm searches for the next shortest path and repeats the same wavelength search procedure as described above. As an effect, the algorithm will first use the wavelengths in the lower-indexed space and leave the wavelengths in the higher-indexed space with higher availability for future requests, thus all of the in-use wavelengths are packed towards the lower end of the wavelength space. Furthermore, since shorter paths are preferred, power consumption and wavelength occupation of the lightpath are reduced.



(a)



(b)

Fig. 4. Inter data center testbed: (a) Configuration, and (b) Picture. The testbed consist of 3 data centers in 5–25 km distance, each equipped with 3 racks and 3 servers. The data center and metro network control planes are running on separate SDN servers.

4. Testbed implementation

We built a three-node data center testbed using commodity optical and electronic components to evaluate the architecture. Figure 4(a) demonstrates the testbed configuration. Each data center is connected to the other two with either 5 or 25 km of Single-Mode Fiber (SMF). Each data center consists of 3 racks, supporting 3 servers. Servers are equipped with a 10G Network Interface Card (NIC), an Intel Xeon 6-core processor and 24 GB of RAM, running Scientific Linux 7 (CentOS 7). For the ToR switches, we used Pica8 OpenFlow Ethernet switches running OpenVSwitch. ToRs in each data center are aggregated in a hybrid architecture using a L2/L3 Ethernet switch (Juniper EX-4500, D-Link DGS-108) and an optical gateway. The optical gateway is implemented using OSS (Calient S320, Polatis 10 and 16), Nistica WSS and 1:8 DWDM

Table 1. End-to-End Delay to Establish a Connection Between Racks of Two Data Centers, Measured on the Testbed.

Component	Delay (ms)
Total Northbound API (Redis)	4.55
Optical Gateway	72
Routing and Wavelength Assignment Algorithm	10
Metro Network Controller Code	1.7
Total	88.25

Mux/Demux supporting C26, C28, C30, C32, C34, C36, C38, and C40 ITU wavelengths. The optical transceivers are DWDM 10G SFP+ ZR modules, providing 24 dB optical link power budget. Figure 4(b) is a picture of the implemented testbed.

For the control plane, each data center has a controller server and there is also one for the metro network. Controllers are implemented in-house using Python. For the communication between the controllers and also for the north-bound API inside data centers, we implemented a pub/sub messaging system using Redis [15]. The southbound APIs of the optical gateway are TL1 commands for the Polatis and Calient optical switches and in-house developed C code for the Nistica WSS. Commercial SDN controllers such as RYU [13] can be used to improve the performance and also integrate the electronic (OpenFlow) and optical switching.

5. Evaluations

In section II, we demonstrated the switching time of the software-defined optical gateway. In this section, first we evaluate the implementation of the control plane by measuring the total control plane overhead to establish a connection. Then we perform end-to-end evaluations on an inter data center testbed by measuring the connections throughput, migrating VMs and bulk data transfer. We conclude this section with simulation results on scalability of the wavelength and routing assignment algorithm.

5.1. Experimental results

We evaluated the control plane implementation by measuring the total delay to establish a new connection. It consist of the northbound API for the application to control plane layer connectivity, implemented using Redis, the optical gateway reconfiguration delay (OSS and WSS switching time that are configured simultaneously), the wavelength and routing kernel, implemented by a *first fit* algorithm and the controller code. Table 1 demonstrates the values that is averaged over 20 measurements. The northbound API delay is measured over campus internet. The total delay is under 100 ms that is mainly due to the optical gateway reconfiguration time. This delay is reasonable for establishing an inter data center connection in the physical layer.

We performed experimental evaluations in the architecture demonstrated in Figs. 3(a) and 4(a). First, we measured the throughput of direct connections between servers of DC₁ and DC₂. Connections are established using DWDM channels C26 and C28. Then, we established connections between the same two servers of DC₁ and DC₂, but indirectly through DC₃. The optical transceivers have nominal throughput of 10 Gbps. Figure 5(a) demonstrates the throughput values over time. All connections either direct or indirect achieved close to full capacity throughput, consistently.

Next, we performed an end-to-end evaluation by requesting VM migrations from DC₁ to DC₂. In the metro network controller, the maximum WDM channel capacity is set to two (C26 and C28). The connection request includes the IP address of the source/destination and the VM name. The VMs configuration is 2 CPU cores, 2 GB of RAM, running Scientific Linux 7. We implemented live migrations by libvirt virtualization APIs. The experiment setup is simultane-

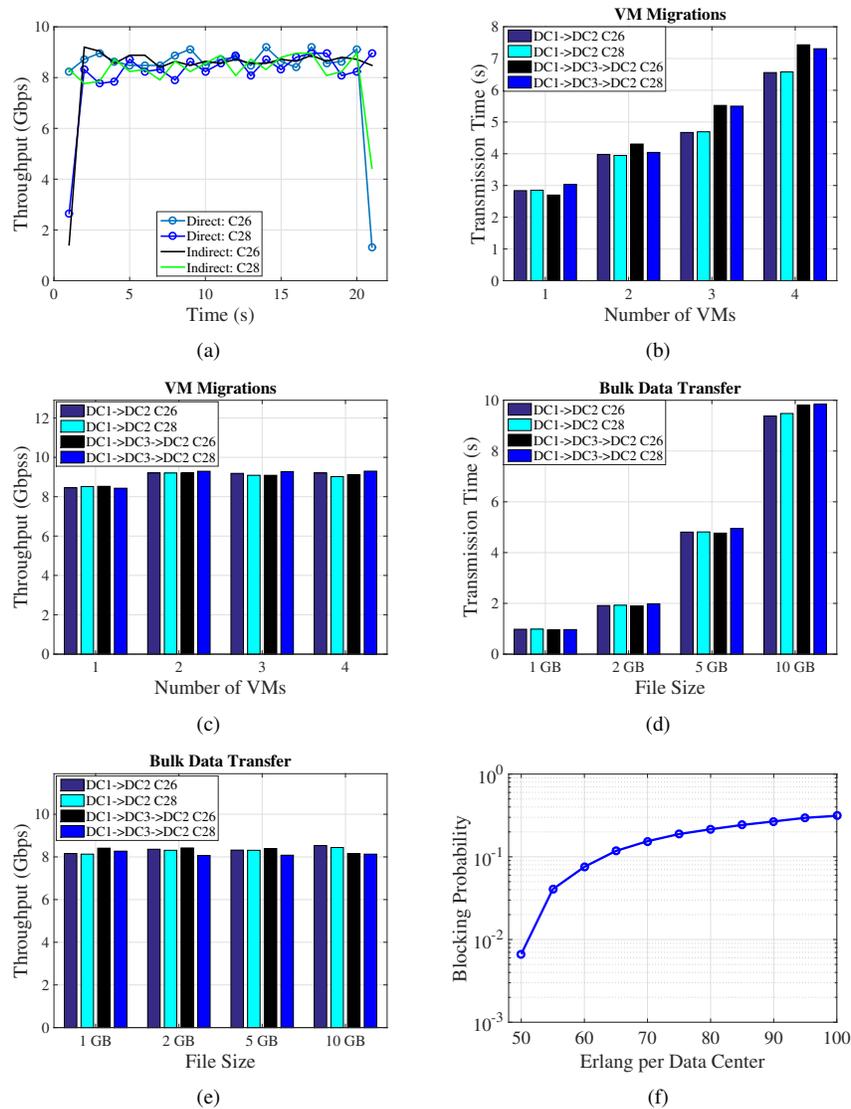


Fig. 5. (a) Measurement of link throughput on the testbed: Establishing 2 direct and 2 indirect (via DC₃) connections between DC₁ and DC₂, WDM channels: C26, C28. (b, c) The transmission time and throughput for 1-10 GB bulk data transfer with direct and indirect connections. (d, e) The transmission time and throughput for 1-4 simultaneous VM migrations with direct and indirect connections. (f) Blocking probability of cross-data center connections in a 4 × 4 mesh network.

ous migration of 1–4 VMs from the server in rack 1 of DC₁ to the server in rack 1 of DC₂. For the indirect experiment, we set the capacity of the link between DC₁ and DC₂ to zero, thus the wavelength and routing kernel finds an alternative route via DC₃. We performed successful live migrations with close to full capability optical link utilization. Figures 5(b) and 5(c) demonstrate the average throughput during migrations and the total migration time. In the case of 2–4 simultaneous migrations, the link is saturated, thus they all have similar throughput. The slight difference in the throughput/transmission time of the wavelengths in the same link is due to

using optical transceivers with slight difference in characteristics.

In the final experiment, we performed bulk data transfer between the servers of rack 1 in DC₁ and DC₂. The setup is similar to the previous experiment (supporting 2 WDM channels per link and setting DC₁–DC₂ link capacity to zero for indirect). The data transfer was implemented using iperf [16] that is a common network performance measurement tool. Figures 5(d) and 5(e) demonstrate the transmission time and the average throughput during the transmission, which confirms successful implementation and close to full capacity optical link utilization.

5.2. Simulation results

In order to evaluate the scalability of the *first fit* algorithm used for the wavelength and routing assignment, we developed a simulation platform. We chose a 5×5 mesh inter data center network as an average optical metro network size and 1000 racks ($n = 1000$) for the size of a small to mid-sized data center in each node of the network. Data centers request rack-to-rack cross data center connections according to a Poisson process with arrival rate p . The destination is determined with uniformly random distribution. Connections have an exponentially distributed holding (service) time with mean h . The simulation hence emulates a birth-death process on the entire network, with an *Erlang* load denoted by $E = nph$.

Figure 5(f) shows the blocking probability of the network (y-axis) under different Erlang loads (x-axis), when each link has 80 available wavelengths. The change in load is obtained by varying the per-rack request rate p while fixed the mean holding time h to 10 (minutes). As the result shows, when the load is less than 50 Erlangs, i.e. requests per data center per minute is less than 5, the network does not see a blocked request that is an acceptable performance [14].

6. Discussion

Nowadays, enterprises and cloud-providers are progressively deploying small and mid-sized data centers. In addition, large-scale data centers are reaching the scalability limits due to the network size and power consumption. Therefore, eventually data centers will also require to scale out in distance and applications need to distribute the operation over several data centers. Optical networks provide the inter data center connectivity with the terabit DWDM capacity. However with the exponential increase in data generation, for efficient operation, a reconfigurable network with an intelligent control plane is necessary.

Flexibility and rapid reconfigurability are the main requirements of the optical network for inter data center connectivity. We tried to address these issues by an architecture capable of bandwidth on-demand and low latency connection setup time. The former metric can be further improved by leveraging adaptive modulation format transceivers. The latter, by using advanced optical switches [17, 18, 19] in the optical gateway, implementing faster APIs and increasing processing power of the controller servers. Faster switching time can improve stability of existing WDM channels in the physical layer by mitigating transient effects of EDFAs. Also, machine learning techniques can be used in the control plane to add intelligence to the network for further optimization and efficiency.

7. Conclusion

Inter data center connectivity is essential for enterprises to get the most advantage of Big Data. We presented a software-defined metro inter data center network that supports transparent and bandwidth selective connections, on-demand. Combining optical circuit switching and WDM provides terabits of bandwidth with wavelength granularity and data rate transparency in switching. It also avoids O/E/O conversions at the switches and improves energy efficiency. Furthermore, the SDN architecture has released the power of optics to enable a fast and reconfigurable interconnection network to improve optical link utilization. This architecture enables

scaling smaller data centers in distance to achieve higher computing power and improving application reliability by distributing the operation.

Acknowledgment

This research was supported in part by NSF ERC for Integrated Access Networks (CIAN) (EEC-0812072) and NSF CNS Networking Technology and Systems (NeTS) (CNS-1423105). We would also like to thank Juniper Networks, Calient Technologies and Polatis for their generous donations to our data center testbed.