

RADWAN: Rate Adaptive Wide Area Network

Rachee Singh
rachee@cs.umass.edu
University of Massachusetts, Amherst

Manya Ghobadi
mgh@microsoft.com
Microsoft Research

Klaus-Tycho Foerster
klaus-tycho.foerster@univie.ac.at
University of Vienna

Mark Filer
mafiler@microsoft.com
Microsoft

Phillipa Gill
phillipa@cs.umass.edu
University of Massachusetts, Amherst

ABSTRACT

Fiber optic cables connecting data centers are an expensive but important resource for large organizations. Their importance has driven a conservative deployment approach, with redundancy and reliability baked in at multiple layers. In this work, we take a more aggressive approach and argue for adapting the capacity of fiber optic links based on their signal-to-noise ratio (SNR). We investigate this idea by analyzing the SNR of over 8,000 links in an optical backbone for a period of three years. We show that the capacity of 64% of 100 Gbps IP links can be augmented by at least 75 Gbps, leading to an overall capacity gain of over 134 Tbps. Moreover, adapting link capacity to a lower rate can prevent up to 25% of link failures. Our analysis shows that using the same links, we get higher capacity, better availability, and 32% lower cost per gigabit per second. To accomplish this, we propose RADWAN, a traffic engineering system that allows optical links to adapt their rate based on the observed SNR to achieve higher throughput and availability while minimizing the churn during capacity reconfigurations. We evaluate RADWAN using a testbed consisting of 1,540 km fiber with 16 amplifiers and attenuators. We then simulate the throughput gains of RADWAN at scale and compare them to the gains of state-of-the-art traffic engineering systems. Our data-driven simulations show that RADWAN improves the overall network throughput by 40% while also improving the average link availability.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

SIGCOMM '18, August 20–25, 2018, Budapest, Hungary

© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-5567-4/18/08...\$15.00

<https://doi.org/10.1145/3230543.3230570>

CCS CONCEPTS

• **Networks** → *Physical links; Traffic engineering algorithms; Network economics; Network performance analysis; Network reliability; Wired access networks;*

KEYWORDS

Traffic Engineering, Wide Area Networks, Optical Backbone

ACM Reference Format:

Rachee Singh, Manya Ghobadi, Klaus-Tycho Foerster, Mark Filer, and Phillipa Gill. 2018. RADWAN: Rate Adaptive Wide Area Network. In *SIGCOMM '18: SIGCOMM 2018, August 20–25, 2018, Budapest, Hungary*. ACM, New York, NY, USA, 14 pages. <https://doi.org/10.1145/3230543.3230570>

1 INTRODUCTION

Optical backbones are million-dollar assets, with fiber comprising their most expensive component. Companies like Google, Microsoft, and Facebook purchase or lease fiber to support wide-area connectivity between distant data center locations but have not been able to fully leverage this investment because of the conservative provisioning of the optical network. We show that wide area fiber links exhibit significantly better signal quality (measured by the signal-to-noise-ratio or SNR) than the minimum required to support transmission at 100 Gbps, leaving money on the table in terms of link capacities.

In other words, there is potential to operate fiber links at higher capacity, thereby increasing the throughput of existing optical networks. We analyze historical SNR from 8,000 optical channels in a backbone network and find that the capacity of 64% of the links can be augmented by 75 Gbps or more, leading to a capacity gain of over 134 Tbps in the network. However, we argue that simply raising link capacities to a higher value (*e.g.*, 150 Gbps or 200 Gbps) increases the rate of link failures because the signal quality fluctuates and operating near the SNR threshold makes links susceptible to failure.

Moreover, enforcing a static link capacity forces operators to treat link failures as *binary* events: when the SNR of a link

falls below a static threshold, the link is treated as “down.” We show this is wasteful, as at least 25% of current failures can be mitigated by reducing the rate of transmission from 100 Gbps to 50 Gbps.

At the core of these issues is a fundamental orthodoxy in the operation of wired networks: a fiber link is either up with a fixed capacity or it is down, largely oblivious to changes in the quality of the underlying optical signal. In this school of thought, operators are forced to account large margins between the actual SNR and the operating capacity if they want to avoid frequent link failures. In contrast, wireless networks employ a variety of schemes to adapt the transmission rate in response to changing signal quality [3, 26, 30]. However, adapting transmission rates to the wireless channel quality is difficult, as the quality can vary at time-scales shorter than a single packet transmission time [30]. In addition, obtaining accurate measurements of received signal strength indication (RSSI) of wireless media (a proxy for true SNR) is difficult in practice [15] because of issues like miscalibration and packet corruption.

We argue that optical links are well positioned to be rate-adaptive. First, signal quality varies at a much coarser time granularity in fiber than in wireless media (hours as opposed to milliseconds). This stability can be leveraged in wide area networks to amortize the cost of infrequently shifting between multiple discrete modulation schemes as signal quality changes. Second, unlike wireless signals, optical signal quality is easily inferred from the bit-error rate (BER) reported after forward error correction (FEC). Leveraging these benefits, we present RADWAN (Rate Adaptive WAN), a system that adapts channel bit-rates in WANs to improve the overall throughput and availability of the network.

RADWAN consists of a centralized rate-adaptive WAN controller that gathers SNR from all fiber channels in the network to adjust the modulation format of the channels to achieve higher or lower data rates. In traditional wide area settings, the QPSK modulation format supports data rates of 100 Gbps for distances up to 3,000 km, 8QAM allows 150 Gbps for distances up to 2,100 km, and 16QAM allows 200 Gbps for distances up to 800 km (see §7 for a discussion on distance). By switching links to a lower modulation format (e.g., BPSK with data rates of 50 Gbps), RADWAN allows critical WAN links to function at lower data rates instead of failing altogether. We refer to these variable capacity links in RADWAN as *dynamic capacity links*. By building on top of existing software-based WAN controllers [16], RADWAN allows traffic engineering schemes to exploit dynamic capacity to improve network throughput. We make two key contributions to make rate adaptive WANs practical:

Optimal WAN traffic engineering. A major challenge associated with dynamically adapting link capacities in WANs is the latency incurred by network hardware when changing a link’s modulation format. To reconcile the latency of capacity changes with the benefits of adapting link capacities in WANs, the RADWAN controller re-formulates the centralized traffic engineering optimization problem to avoid unnecessary capacity reconfiguration (§4). We evaluate our controller by comparing the throughput gains of employing RADWAN at scale to those of a state-of-the-art controller. Our results show that in a real-world network topology and with conservative traffic churn settings, RADWAN improves the overall network throughput by 40% while also improving the link availability (§6). We estimate that RADWAN lowers the dollar per gigabit per second cost of traffic by 32% (§7).

Avoiding high latency of modulation reconfiguration. We build a testbed emulating a WAN connecting four data centers via 1,540 km of fiber. Using this testbed, we confirm the viability of modulation reconfiguration to achieve greater network throughput. We benchmark the behavior of the RADWAN controller as it reacts to SNR degradation by switching to a lower modulation format. During the modulation change, the line-rate traffic on the affected link is migrated to a backup path until the modulation change is complete (§5). Our experiments show that reconfiguring modulation formats on commodity hardware incurs a latency of 68 seconds, on average. We develop a prototype that demonstrates the feasibility of decreasing this reconfiguration time by a factor of 1,000 (§7.1).

RADWAN opens the door to revisiting several classical networking problems in light of dynamic capacity links. For instance, are there graph abstractions that can capture networks with dynamic capacity links? How do classical networking algorithms (such as the maximum-flow problem [11]) change in the presence of variable link capacities? Are there smart capacity planning, failure-recovery, load-balancing, or on-demand bandwidth allocation algorithms that can benefit from rate adaptive links? RADWAN prepares the ground for thinking about these problems.

2 QUANTIFYING THE OPPORTUNITY

We investigate the signal quality of 8,000 optical channels in a large optical backbone network. Our dataset consists of the average, minimum, and maximum SNR per channel, aggregated over 15 minute intervals for \approx three years (Feb. 2015 - Dec. 2017). We characterize the SNR of these channels and quantify its variations. In wireless networks, signal quality may vary in short time intervals and estimating SNR is complicated by signal interference [30], but signals in fiber optical media do not face these challenges.

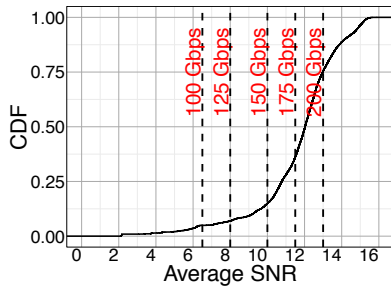


Figure 1: Distribution of the average SNR of over 8,000 channels in a backbone network for three years. Note that the SNR of optical channels is much higher than the required SNR for a 100 Gbps bit rate (6.5 dB).

Impact on capacity. We note that in our data, the average SNR is much higher than the required threshold for operating links at their current rate of 100 Gbps. Figure 1 shows the distribution of the average SNR of all 8,000 channels, with vertical dashed lines marking the SNR threshold for various rates. The figure shows that 95% of the channels have an average SNR above the required threshold for 100 Gbps. Even better, 64% of the channels have an SNR that can support data rates of 175 Gbps or higher but are currently used for a conservative 100 Gbps only. This represents a significant opportunity to improve the throughput of optical links by operating closer to the actual SNR of the signal.

But what about stability? While the average SNR may be well above what is needed to drive links at 100 Gbps, simply increasing the links' rate to a higher value, say 150 Gbps, will not work in practice because SNR fluctuates, as illustrated in Figure 2. The figure shows the SNR of 40 channels on one fiber cable observed over 2.5 years. We note that the SNR of these channels is largely stable, but there are occasional dips caused by impairments in the fiber or other optical hardware. The frequency and duration of these dips vary for different fibers in the network. The dashed horizontal lines in Figure 2 show the required threshold for various data rates. Higher SNR means we can have higher data rates.

We further consider the variability of SNR across all links for different time-scales. For each time interval of size 15 minutes, 10 hours, 1 day and 1 week, we calculate the variability of SNR (the difference between maximum SNR and minimum SNR) for all optical channels in the backbone network. Figure 3 shows the distribution of SNR variation in time intervals of different sizes. We confirm that SNR remains stable over several hours at a time. A small fraction ($\leq 5\%$) of links show a variation of over 1 dB in the 10 hour interval. Moreover, although our SNR measurements are aggregated over 15 minute intervals, we argue that our conclusions are sound, as Figure 3 shows negligible variation

in SNR for 15 minute time intervals. This contrasts with wireless media where significant SNR changes can happen within a few milliseconds.

Why do we need variable bandwidth links? Based on our observation of the mostly stable but over-provisioned SNR of links, one might be tempted to operate links closer to the actual SNR by simply making a one-time decision to increase the transmission rate of all links. However, we find that the frequency of link failures increases if we cannot dynamically adapt to SNR changes. This is because infrequent but sizable variations in SNR occur in fiber links. While the SNR of a small fraction of links changes significantly in a few hours, 10% of all links undergo 2 dB of change in SNR within a week (Figure 3). To illustrate this, we select a fiber where the SNR of each link (i.e., optical channel) is high enough to make all capacity denominations feasible over three years. We then analyze the number of failures the links would undergo if they were modulated with higher but static capacities. Figure 4(a) shows that links on this fiber do not see a significant increase in the number of failures as the capacity is increased up to 175 Gbps, but some would have up to 100 failures if driven at 200 Gbps. We find this behavior repeated in other fibers, but depending on the number of links, fiber length, technology, and age of equipment, the point at which the failures start to increase differs for each fiber and for each channel on the fiber. Hence, it is impossible to select a one-size-fits-all static capacity that is higher than 100 Gbps.

Next, we characterize the duration of SNR dips to evaluate the magnitude of disruption they could cause if we choose a higher modulation (hence higher bandwidth). Figure 4(b) plots the duration of link failures for the various modulated bandwidths (based on the link's average SNR). We observe that such SNR dips last for several hours which means we cannot simply select a static modulation and dismiss the SNR dip events. The good news is that by enabling variable bandwidth links, we can react to SNR dips by changing the bandwidth to match the SNR.

Impact on availability. Today, when the SNR of a link's optical signal drops below its pre-determined threshold, the link is declared down. However, not all failures are complete loss-of-light. SNR drops may be caused by planned maintenance work (e.g., a line card replacement) or unplanned events (e.g., fiber cut, hardware failure, human error). While some of these impairments make the link unusable (e.g. fiber cuts), others may simply lower the signal quality (e.g. degradation of an amplifier) without completely shutting off the signal. Links undergoing failures due to lowered signal quality can still be used to send traffic at a reduced rate, highlighting another opportunity to improve link availability.

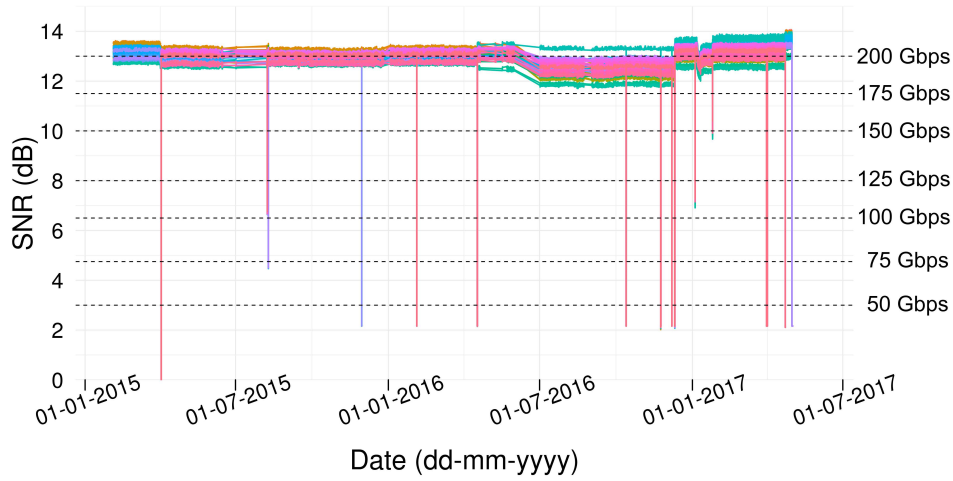


Figure 2: SNR variations in 40 optical channels (i.e., IP links) on a wide area fiber cable. Dotted lines represent the feasible link capacity for a particular SNR.

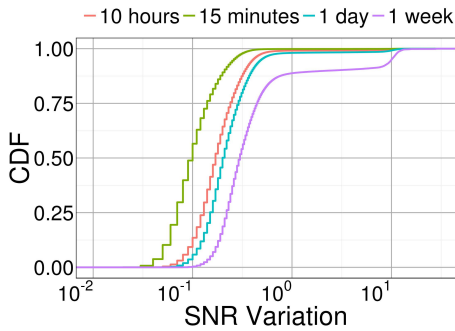


Figure 3: Variations in the channel SNR in intervals of different durations. Observe that most links do not observe significant variation in SNR for several hours.

To define the opportunity area, we record the lowest SNR of failure events (when the SNR falls below the 100 Gbps threshold which is 6.5 dB). Figure 4(c) shows the distribution of lowest SNR values at link failures. We observe that in 25% of the failures, the lowest SNR is above 3 dB, enough to drive a link at 50 Gbps capacity. Therefore, 25% of the link failures could have been avoided by driving the impacted links at 50 Gbps, indicating the improvement in availability offered by dynamic capacity links.

3 DYNAMIC CAPACITY LINKS

Our characterization of the SNR values of optical links suggests they are currently operating well below their potential transmission rates. However, operating links at constant transmission rates closer to the observed SNR increases the likelihood of link failures. To balance this trade-off, we propose a dynamic adjustment of physical link capacities in *centrally*

controlled wide area networks by changing the *modulation format* of optical signals. These choices are motivated by the latest hardware and software developments in the industry:

Adapting bit-rates by changing the modulation. Recent advances in the development of bandwidth variable transceivers (BVTs) provide a promising first step towards increasing the network bandwidth and improving availability by decreasing transmission rates in the face of low SNR (vs. incurring a link failure). State-of-the-art BVTs are capable of modulating signals on the fiber with three different formats: 16QAM, 8QAM and QPSK. All other factors being constant, signals in 16QAM format can carry traffic at 200 Gbps, 8QAM can carry 150 Gbps and QPSK can carry 100 Gbps. However, these transceivers were designed with the assumption that operators would make a one-time choice of modulation format.

This is reflected in the latency incurred in changing the modulation of ports on modern Arista 7504 switches. In our experiments, we find that on average, changing the modulation of a port incurs a latency of over one minute. During this time, the link undergoing the modulation change is down and cannot carry traffic. This is because of the assumption by the manufacturers that the modulation change is a one-time event. To benchmark the reconfiguration latency, we experiment with a transceiver evaluation board and investigate ways of reducing capacity reconfiguration time (Section 7.1). We note that it will take significant engineering efforts to make hitless capacity change production ready for use.

Software Driven WANs. Effective utilization of network infrastructure in modern WANs is enabled by software-driven

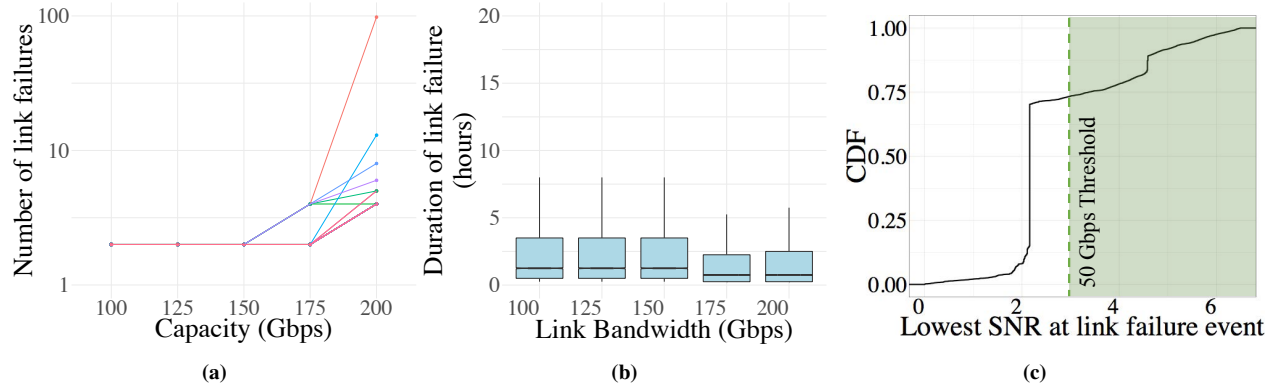


Figure 4: (a) Number of link failures for 40 links (one color per link) for a given capacity. For this particular fiber, while increasing capacity up to 175 Gbps does not increase link failure events, achieving 200 Gbps capacity comes at the cost of increased link failures. (b) Duration of failures if WAN links operate at a given capacity. (c) Distribution of the lowest SNR values when a link failure event happens. The lowest SNR is above 3.0dB (sufficient to drive a link at 50 Gbps) 25% of the time.

centralized traffic engineering (TE) [16, 17, 23] that maximizes the network throughput for changing demand matrices. Therefore, we consider the implementation of dynamic capacity links in such networks. We note that throughput maximization is one possible goal of traffic engineering. Previous work has formulated TE to achieve optimal social welfare [18], meet deadline-sensitive transfers [22], and provide improved guarantees for high priority traffic [16].

TE controllers are consumers of network link capacities, as they make decisions about routing flows along the best paths with available capacity. Had link capacity reconfiguration been a hitless phenomenon, existing TE controllers could largely function unmodified with dynamic capacity links. However, capacity reconfiguration is expensive, as it causes a link outage lasting for over a minute. We discuss the impact of this additional constraint on TE controllers in the next section.

4 TRAFFIC ENGINEERING WITH DYNAMIC CAPACITY LINKS

In a network with dynamic capacity links, the state of the network in each run of the TE optimization algorithm is dependent on the links' underlying SNR. Therefore, TE controllers must be modified to gather the SNR of all links in the network and to treat the link capacities as variables. Our proposed RADWAN centralized TE controller can leverage dynamic capacity links to achieve higher network throughput and availability. RADWAN handles a spike in the demand matrix by upgrading the capacities of one or more links. However, state-of-the-art bandwidth variable transceivers (BVTs) require over a minute to change the capacity of a link (§ 7.1), rendering the link unusable for that period. In response to this link flap, existing traffic flows must be

migrated away from the link undergoing capacity reconfiguration, but such flow migrations can cause transient congestion in the network and must be done minimally.

Therefore, we argue that in a network composed of dynamic capacity links, the objective of traffic engineering changes from simply maximizing the network throughput to maximizing throughput *while minimizing churn caused by link capacity reconfigurations*. In Section 4.1 we discuss how network churn can be quantified to achieve low disruption while meeting traffic demands via link capacity reconfiguration.

4.1 Quantifying network churn

Current hardware does not support hitless capacity changes; therefore, we propose dealing with churn induced by link capacity changes in software. As a first step, we introduce a definition of churn induced by a link capacity change in terms of the rate of traffic on the link. The capacity change (either an increase to meet demands or a decrease due to lowered signal quality) of link l carrying f_l units of traffic will displace f_l units. The displacement of large flows is more likely to cause transient congestion as opposed to smaller flows. Therefore, we define churn induced by the capacity change of link l as:

$$churn(l) = f_l \quad (1)$$

The overall churn induced by capacity changes in a network, C , is a summation of the churn from each link undergoing capacity change:

$$C = \sum_{links} churn(l) \quad (2)$$

We note that this is only one of many possible ways to define the churn caused by link flaps in the network. We encourage practitioners to consider other definitions to reduce

the churn of preferred traffic classes (e.g., interactive traffic over background traffic).

4.2 Computing flow allocations

When computing allocations of flows along different paths in a network composed of dynamic capacity links, the goal of RADWAN is to maximize the network utilization (as was the case with earlier work [16, 17, 23]) while keeping churn due to capacity reconfigurations minimal. In this section, we formulate this goal as a constrained optimization problem using the definition of churn from Section 4.1. RADWAN periodically evaluates the optimization goal to assign traffic flows along network paths. In each round of its operation, RADWAN has access to attributes of the network state which serve as input to the optimization problem. We now describe various elements of the RADWAN controller.

Inputs. Traditional TE controllers take *network topology* and *traffic demand matrix* as input to compute allocations of flows along label-switched network paths. In addition to these, our controller requires SNR measurements for all physical links in the network. Using this information, the controller derives the *potential* capacity of each link, over its existing capacity.¹ Implicitly, the controller is also aware of the existing flow on all links in the network, assigned in the previous round of controller operation.

Allocation Objective. Algorithm 1 describes the optimization goal of RADWAN. At its core, the optimization is a modified multi-commodity flow that maximizes overall throughput of the network while augmenting link capacities minimally. The optimization variables $b_{i,j}$ specify the allocation of flow i along path j in the network. Allocation of flow along link l in the network is constrained by the sum of the link capacity (c_l) and the potential increase in capacity (p_l) depending on the link's SNR. ϵ is a small positive constant denoting the relative importance of the two aspects of the objective function: maximizing throughput and minimizing churn. Finally, in a given round, the network churn caused by the capacity change of link l is 0 if the optimal flow assigned to the link is less than or equal to the link's capacity (c_l). However, if the link has more flow assigned to it than its current capacity, it induces network churn equal to the amount of traffic on it (f_l), as assigned in the previous round of flow allocation. The nature of network churn makes the objective function of the optimization piece-wise linear.

Approximation to Linear Program. To efficiently solve the optimization objective described in Algorithm 1, we approximate the definition of churn as:

¹Even if there is potential to increase a link's capacity by, say, 50 Gbps, the controller must do an upgrade only if this extra capacity is needed to meet traffic demands.

Algorithm 1: Traffic Engineering Optimization

1 Inputs:

- 2 d_i : flow demands for source destination pair i
- 3 c_l : capacity of each link l
- 4 p_l : potential capacity increase of each link l
- 5 $I_{j,l}$: 1 if tunnel j uses link l and 0 otherwise
- 6 f_l : existing flow on link l ($f_l \leq c_l$)
- 7 T_i : set of tunnels set up for flow i

8 Outputs:

- 9 $b_i = \sum_j b_{i,j}$: b_i is allocation to flow i
- 10 $b_{i,j}$ is allocation to flow i along tunnel j

11 Maximize: $\sum_i b_i - \epsilon(\sum_l \text{churn}(l))$

12 subject to:

- 13 $\forall i, 0 \leq b_i \leq d_i$
 - 14 $\forall i, j, b_{i,j} \geq 0$
 - 15 $\forall l, \sum_{i,j} I(j,l) b_{i,j} \leq c_l + p_l$
 - 16 $\forall i, \sum_{j \in T_i} b_{i,j} \geq b_i$
 - 17 $\text{churn}(l) = \begin{cases} 0, & \sum_{i,j} b_{i,j} I(j,l) \leq c_l \\ f_l, & \text{otherwise} \end{cases}$
-

$$\text{churn}(l) = \max(0, (\sum_{i,j} b_{i,j} I(j,l) - f_l)) \quad (3)$$

This monotonically increasing value of churn, depending on the flow assignments $b_{i,j}$, is different from the actual churn value which is essentially a step function; however, this reasonable approximation allows us to convert Algorithm 1 to an efficiently solvable linear program.

Managing Churn. For the duration of a link flap, no traffic can be routed along this link. As the impacted links will be offline for just one minute, the affected traffic (churn) has to be managed efficiently to ensure low disruption. We thus compute a single intermediate flow allocation, where the churn is distributed along routes without link flaps. We show in Section 6 that a single intermediate step suffices, as the number of link flaps per reconfiguration is low in practice (see Figure 10). Methods for networks with highly unstable SNR are described in Section 7.1. Once hitless capacity reconfiguration is production ready, the intermediate flow allocation step described in this paragraph can be omitted.

Importance of ϵ . The ϵ parameter defines the balance between RADWAN controller's tendency to maximize network utilization and minimize network churn due to capacity reconfigurations. We encourage operators to use a value of ϵ that captures their willingness towards capacity reconfigurations. We note that future optical equipment that offers reduced capacity reconfiguration time will make capacity changes more attractive and operators can use

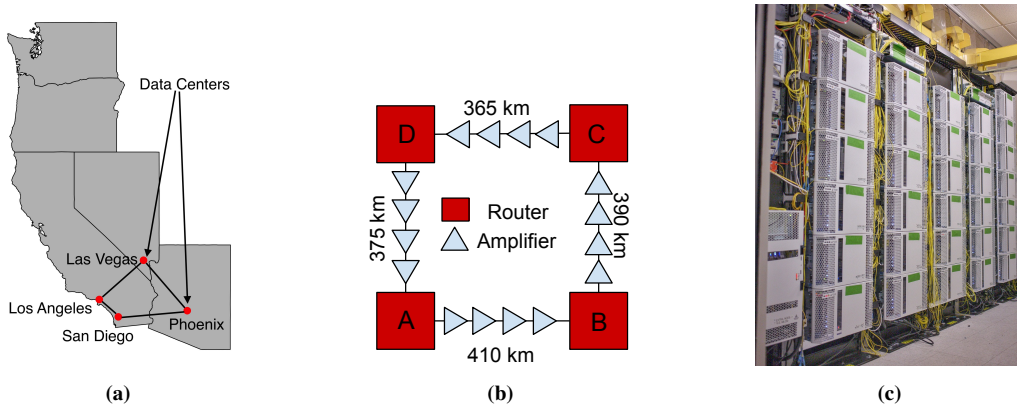


Figure 5: (a) Geographic scale of the testbed built to demonstrate the operation of RADWAN. Our testbed emulates a WAN connecting four major cities on the west coast of the United States. (b) Logical view of the testbed where four routers (logically split from a modular chassis switch) emulate four data centers. These routers are connected via hundreds of kilometers of optical fiber and regularly spaced amplifiers. (c) Photograph of the electrical and optical equipment in our testbed.

smaller ϵ values in the optimization to reflect increased willingness for capacity reconfiguration.

4.3 Controller Implementation

We implement RADWAN, the traffic engineering controller based on the goals outlined in the previous subsection. The controller implements Algorithm 1 using the popular optimization library CVXPY [5] in Python 2.7.

RADWAN computes flow allocations for the input demand matrix in each round of its operation. Before solving the optimization, RADWAN uses the link-level SNR information to determine: (i) links for which the total capacity must now be reduced since the new SNR is too low to support the existing capacity. These *capacity downgrades* must be performed even though they will cause the impacted links to be down for roughly a minute; (ii) the potential capacity of other links, above their current capacity, depending on the SNR of the link. For instance, a link could be operating at 100 Gbps, but if it has an SNR of 10.2 dB, it has a potential capacity increase of 50 Gbps, as its capacity can be augmented to 150 Gbps.

In what follows, we present the results of an extensive testbed evaluation of RADWAN (§5), benchmarking the effect of optically changing links' capacities on the IP layer. We then simulate RADWAN and compare it with our implementation of the SWAN controller as described in [16]. We perform a data-driven evaluation of the behavior and performance of these two controllers in Section 6 and show the gains of capacity variable links on the overall network throughput.

5 TESTBED EVALUATION

In this section, we build a testbed consisting of 1,540 km of fiber and 16 optical amplifiers to evaluate the feasibility of deploying RADWAN in a moderate sized WAN. Our goal is to highlight the impact of modulation changes on realistic traffic flows. We also provide insights to both researchers and practitioners into the state-of-the-art hardware components required to realize a rate-adaptive wide area network.

5.1 Testbed Implementation Details

We build a moderate sized testbed which emulates a WAN interconnecting four data centers, as shown in Figure 5(a), to evaluate RADWAN. Each data center consists of a router connected to its neighbors through hundreds of kilometers of optical fiber. To prevent signal deterioration, we connect Erbium Doped Fiber Amplifiers (EDFAs) at approximately every 65-120 kilometers of fiber length. For simplicity, Figure 5(b) represents the logical view of the WAN.

Note that we had access to only one Arista 7504 modular chassis; therefore, we used Virtual Routing and Forwarding (VRF) [6] to logically split the same physical switch into four routers (named A, B, C and D in Figure 5(b)). Each VRF has a separate routing table and routing protocol instances. By configuring relevant physical interfaces to be in separate VRFs and connecting the interfaces via optical components (fiber, amplifiers), we achieve a logical topology whereby traffic between ports on the switch is sent out on the wire. We verify bi-direction connectivity between each pair of nodes A, B, C and D. The Arista 7504 has integrated bandwidth variable transceivers manufactured by Acacia Inc. (the BVT module, AC 400, is described in detail in Section 7.1). These allow us to configure three modulation formats (QPSK, 8QAM and

16QAM) on the switch ports. The complete testbed, including optical and electrical equipment is shown in Figure 5(c).

We implement the part of the RADWAN controller responsible for configuring the switch using Arista’s PyEAPI [2] framework. With this, we can programmatically configure the modulation formats of different ports, program routes and query status of our commands.

To generate line rate traffic flows in the topology, we use a Spirent traffic generator [28]. With the help of the Spirent device, we program 400 Gbps of TCP traffic flows to test the dynamic capacity links of the testbed.

5.2 Benchmarking the WAN testbed

Reacting to SNR degradation. Optical signals in fiber can become attenuated because of ill-functioning amplifiers, disturbances caused during maintenance windows or even ambient temperature conditions. RADWAN reacts to signal attenuation by switching to a lower order modulation format that can be supported by the degraded SNR. In the laboratory setting, we use a Variable Optical Attenuator (VOA) device to add configurable amounts of noise (measured in dB) so that we can demonstrate signal attenuation. We connect the VOA between routers *A* and *B* in the test topology. On the underlying switch, this connection is implemented by connecting *Ethernet4/1/1* to *Ethernet3/1/1* with 410km of optical fiber. The Ethernet ports are in separate VRFs (not directly connected), so we set up static routing such that traffic sent from one to the other is sent over the fiber connection. Every five seconds, we increase the noise from the VOA by 1 dBm.

We measure the SNR of the signal on each end of the connection and observe that the SNR of the received signal on *Ethernet3/1/1* steadily deteriorates as the level of noise increases (Figure 6). Once the added noise reaches 16 dBm, the transceiver can no longer recover from the increased errors,² and the port goes down. At this point, the controller reduces the modulation format of the port from 16QAM to 8QAM. The modulation change takes approximately 70 seconds to complete. We then resume incrementing the noise level using the VOA. When the noise level reaches 18 dBm, the transceiver can no longer recover from the errors to support 8QAM format, and the port goes down again. Our controller reacts by reducing the modulation format yet again, this time from 8QAM to QPSK. After roughly 70 seconds of down time, the ports come back up with QPSK modulation format. The addition of noise of 23dBm or more renders the link unusable, even in the lowest supported modulation format. At this point, the link has failed, and the failure is irrecoverable with the current set of hardware.

²Acacia BVTs have 15% soft decision FEC enabled by default.

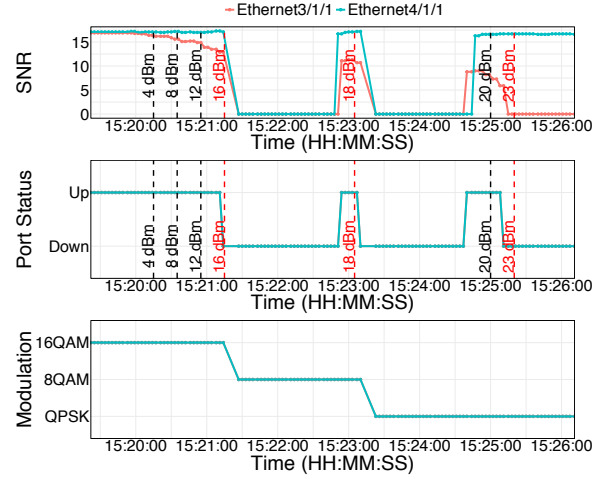


Figure 6: Impact of attenuation on link SNR, port status and modulation format as the amount of signal attenuation increases (shown by dotted vertical lines).

Modulation Change Latency. In the above benchmarking experiment, we changed the modulation format of a link in the testbed in response to SNR degradation. We observe that each change in modulation format changes the status of the ports involved to *down*, making them unavailable for sending and receiving traffic. In Figure 6, we observe that modulation change operations take approximately 70 seconds. This aspect of the latency of modulation change guides the design of the RADWAN controller (Section 4).

5.3 Evaluating Modulation Change

In this section, we demonstrate the capability of RADWAN to react to SNR degradation by reducing the modulation format of ports, allowing links with reduced signal quality to function at lower rates. We provide an end-to-end evaluation of RADWAN as it attempts to meet changing demand matrices by upgrading the capacities of links in the WAN. Additionally, we show that RADWAN migrates flows from a link undergoing capacity up-/downgrade (due to improved/poor SNR) to alternate paths until the modulation change is complete.

In each of the following experiments, we show the transmission rate (Tx Rate) of the traffic we attempt to send between nodes in the topology. An overwhelmed node responds to high traffic volume by dropping a portion of the flows. We capture the net traffic received by the sink node of a flow as the receive rate (Rx rate). In the ideal case, the Tx and Rx rates should match, implying that all the traffic sent by the source is reaching the sink node.

Link capacity upgrade. Figure 7(a) shows the starting state of a network with two flows of 100 Gbps, one from Node *B*

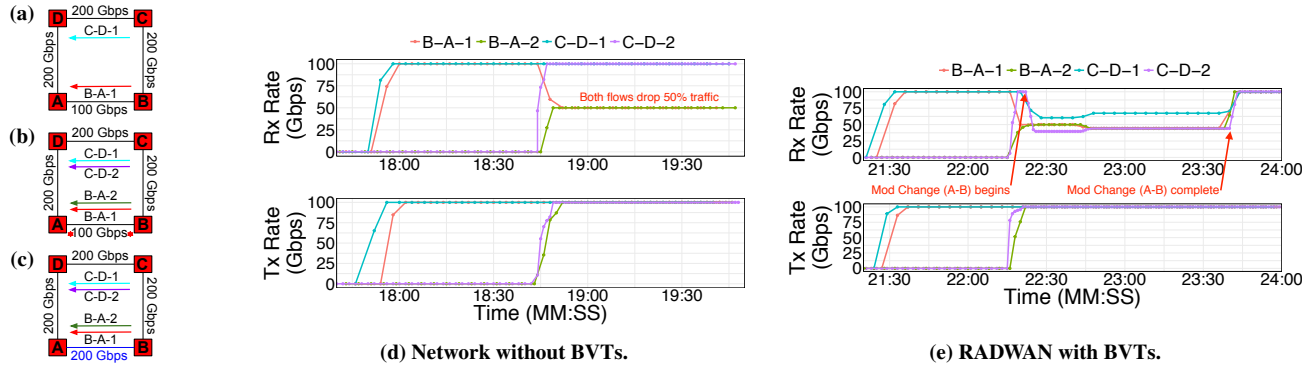


Figure 7: (a) shows the network and link capacities. At the start, all links except link $A-B$ are in 16QAM modulation format, capable of carrying 200 Gbps. $A-B$ being in QPSK format can carry 100 Gbps. In the beginning, there are two flows in the network, each 100 Gbps from $B \rightarrow A$ and $C \rightarrow D$. With an additional demand of 100 Gbps ($B-A-2$) and ($C-D-2$) described in (b), the link $A-B$ gets congested, leading to 50% traffic drops in flows $B-A-1$, $B-A-2$ in the absence of RADWAN, as seen in the Rx Rate in (d). However, in RADWAN deployment, the controller reacts to the increased demand by increasing the capacity of $A-B$ link to 200 Gbps (seen in (c)) by changing the modulation format to 16QAM. While this causes temporary disruption due to rerouting of $B-A$ flows along the $C-D$ link, once the modulation change is complete, the network can carry the flows of 400Gbps without any drops, as seen in the Rx Rate of (e).

to Node A (flow $B-A-1$) and the other from Node C to Node D (flow $C-D-1$). With the introduction of two additional 100 Gbps flows ($B-A-2$) and ($C-D-2$), as shown in Figure 7(b), the network becomes congested, because link $A-B$ can only carry 100 Gbps of traffic. As seen in the Rx rate in Figure 7(d), both $B-A-1$ and $B-A-2$ share the $A-B$ link fairly and drop 50% of their traffic. However, RADWAN can salvage this congestion by increasing the capacity of the $A-B$ link (as seen in Figure 7(c)). To do this, the RADWAN controller reacts to the increased demand by changing the modulation of the $A-B$ link, causing it to be down for roughly one minute. This temporarily congests the $C-D$ link (the Rx rate of all flows drops in Figure 7(e)), because the $B-A$ flows are rerouted. However, once the modulation change is complete, all flows can be transmitted successfully with no packet drops. We note that without augmenting the capacity of link $A-B$, the network could not satisfy 400 Gbps of demand but dynamic capacity links with RADWAN enable us to meet the increased demand.

Link capacity downgrade. Figure 8(a) shows the starting state of our testbed when the network is carrying three flows of 100 Gbps, two from Node C to D ($C-D-1$, $C-D-2$) and one from Node B to A ($B-A-2$). All links in the network can carry 200 Gbps of traffic. Observe that the Rx rate in Figure 8(d) matches the Tx rate, implying there is no packet loss. Now, we attenuate the signal between Node A and B using a VOA device, such that the switch ports can no longer sustain transmission at 200 Gbps. Therefore, the link goes down (Figure 8(b)), causing ($B-A-2$) to be routed over the longer path $B \rightarrow C \rightarrow D \rightarrow A$ which is configured as the backup route. This transition of the $B-A-2$ flow along the longer

path is visible in the utilization of links in the network (Figure 8(d)). Links $B-C$ and $D-A$ are now carrying 100 Gbps of the $B-A-2$ flow (and, thus, are 50% utilized). Note that this leads to congestion on link $C-D$ which can only carry 200 Gbps of traffic; accordingly, it drops 100 Gbps of traffic from the $C-D$ flows. The RADWAN controller can mitigate this congestion by reducing the modulation format of the $A-B$ link to QPSK from 16QAM. It takes roughly one minute for the modulation change to take effect, as observed in the *down* status of link $A-B$ in Figure 8(d). Once the modulation change is complete, link $A-B$ is back up and carries the $B-A-2$ flow without any congestion in the network (Tx/Rx rates match again). The new network state is shown in Figure 8(c). Therefore, our experiments show that RADWAN can react to traffic demands and signal quality by adapting the capacity of links in the WAN.

6 LARGE SCALE EVALUATION

In Section 2, we used three years of SNR measurements to demonstrate that an overall *capacity gain* of 67% is possible by augmenting the capacity of links from 100 Gbps to 125, 150, 175, or 200 Gbps, depending on their average SNR. This is the upper bound of the *throughput gain* achievable with RADWAN. The actual network throughput depends not only on the network state (topology, link capacities, tunnels *etc.*) but also on the traffic demand and acceptable churn (defined in §4). In this section, we simulate the operation of RADWAN in a large backbone network with periodically varying traffic demands to compute the network throughput achieved. We compare the throughput and availability of the network under RADWAN and a state-of-the-art SWAN controller.

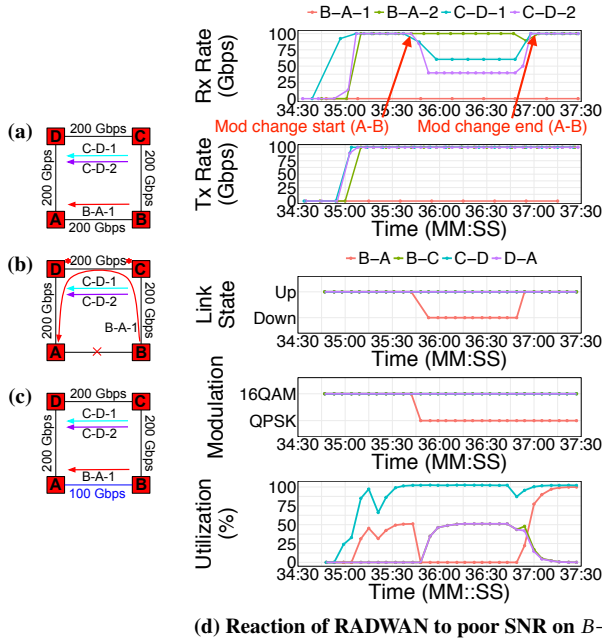


Figure 8: (a) describes the starting network state and link capacities. At the start, all links are in 16QAM modulation format, capable of carrying 200 Gbps. There are three flows in the network, each 100 Gbps, one from $B \rightarrow A$ and two from $C \rightarrow D$. Due to signal attenuation, the link $A-B$ fails as seen in (b), causing the $B-A-2$ flow to be routed over the longer path $B \rightarrow C \rightarrow D \rightarrow A$. Observe that the utilization of links $B-C$ and $D-A$ increases in (d). This causes link $C-D$ to become congested and it drops parts of the $C-D-1$ and $C-D-2$ flows (Tx rate falls below the Rx rate in (d)). RADWAN reacts to this situation by reducing the modulation format of link $A-B$ to QPSK as is allowed by the lowered SNR of the link (see (c)). Once the modulation change is complete, all flows are routed along direct paths without any packet loss, as confirmed by the Tx/Rx rates and link utilizations in (d).

Both controllers are aware of the underlying signal quality of links. But unlike SWAN, RADWAN uses the SNR to update link capacities, choosing amongst discrete choices of 50, 100, 125, 150, 175 and 200 Gbps. As outlined in the previous section, RADWAN only upgrades the capacity of a link to meet increased traffic demand that cannot be met otherwise. Capacity downgrades are done to prevent link failures such that the lower quality link can continue to function at a reduced rate.

6.1 Simulation Setup

We consider the network topology of a large commercial WAN and gather SNR measurements from the optical fiber connecting the nodes in the topology for four randomly chosen days in 2016 and 2017. Both RADWAN and SWAN compute flow allocations along various network paths to meet an elastic demand between each pair of nodes in the network.

Since our WAN currently operates links at 100 Gbps, we consider the performance of SWAN in a fixed capacity network where each link operates at 100 Gbps if the SNR is above the threshold of 100 Gbps modulation; otherwise, the link is down. We refer to this scheme as SWAN-100 in the analysis. However, operators can be more aggressive by operating links at a fixed but higher capacity of 150 Gbps. We refer to SWAN operating in such a network as SWAN-150. SWAN-150 is used to compare the benefit of using rate adaptive schemes like RADWAN over a network with higher but fixed link capacities. While current hardware limitations prevent hitless capacity changes, we simulate the performance of RADWAN under both hitless (RADWAN-HITLESS) and non-hitless (RADWAN) link capacity change behavior.

The traffic demand between each node pair varies periodically every two minutes (demand pattern shown in Figure 9(a)). Our choice of network demands is similar to previous work [16], since rapid changes in demand matrices stress test the TE controllers. We also offset the traffic demand between each pair of nodes by using a randomized value to ensure that at any given point in time, there is sufficient variety of demands in the network.

Simulation Parameters. Unless otherwise stated, the control loop of both controllers is executed every 30 seconds as stated in [16]. In addition, we assume the demand between each pair of nodes can be split across $k = 2$ shortest paths between the nodes. For RADWAN, we set the churn trade-off parameter ϵ (defined in §4) to a conservative value of 0.001. We perform several runs of this experiment, with each run lasting for one day. We find that across four randomly chosen days, our results are similar. Hence, for the sake of brevity, the figures show results from one experimental run.

6.2 Evaluation Metrics

We focus on the following three key aspects of cost-efficient network design to evaluate RADWAN.

Network Throughput. First, we compute the optimal network flow that RADWAN can achieve in each run of the controller and compare it with the optimal flow that SWAN achieves for the same network conditions. This provides the network throughput enabled by both controllers for each run of their control loops for the duration of a day. Figure 9(b) shows the network flow for both RADWAN and SWAN for two hours of a day (zooming into two consecutive hours, picked randomly for the sake of better visibility in the figure). We observe that RADWAN manages to push 40% more traffic than SWAN-100 in the same network. The same observation holds consistently with other hours and days we simulated.

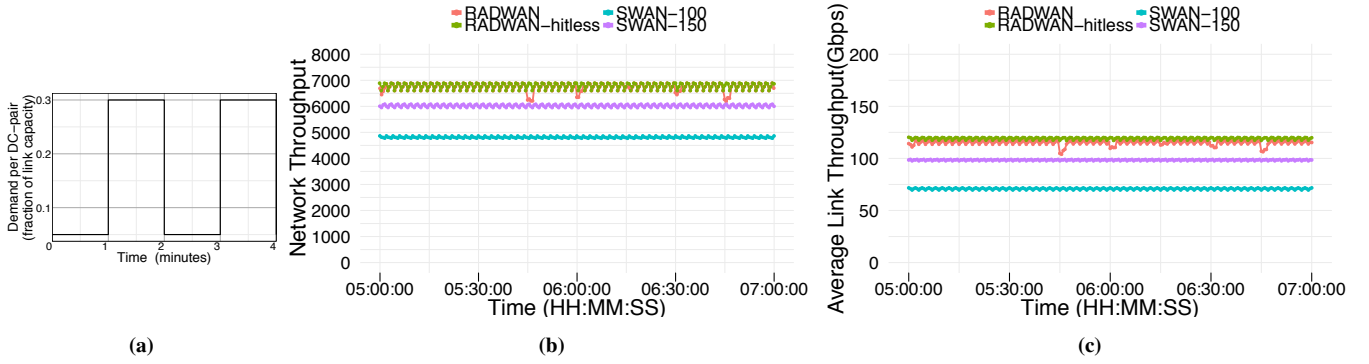


Figure 9: (a) Traffic demand pattern between each pair of nodes in the network similar to prior work [16]. (b) Optimal network flow achieved by different traffic engineering schemes. For better visibility, (b) zooms into two hours of the simulation period. RADWAN achieves 40% higher network throughput than the state-of-the-art mechanism, SWAN-100 (RADWAN and RADWAN-hitless are overlapping curves on top of the graph). We also compare SWAN’s performance with fixed capacity links operating with a static 150 Gbps modulation format. While SWAN-150 provides an improvement over SWAN-100, RADWAN achieves 12% higher throughput than SWAN-150. (c) Average per link throughput. We observe RADWAN achieves 68% higher per link throughput than SWAN-100.

Link Throughput. Next, we compare RADWAN and SWAN’s per link throughput. For each run of the TE control loop, we compute the total traffic carried by each link and average it over all links in the network. Figure 9(c) shows the distribution of average link throughput over time (zoomed over two hours for better visibility). We find that, on average, RADWAN increases the utilization of network links by 68% compared to SWAN-100, getting more utility from each link in the network.

Link Availability. We compute the downtime of links in the WAN as the fraction of total simulation time for which a link is unavailable to carry traffic. Since the WAN we analyze is production grade, it was highly available during the 4 randomly chosen days in this simulation. Therefore, even under the existing SWAN-100 scheme the average link downtime is very small. However, we find that RADWAN reduces the average link downtime by a factor of 18 when compared to SWAN-100 operating in the same network. This is because RADWAN adapts links to lower capacities, when possible, instead of failing them when the signal quality degrades. Even though RADWAN’s capacity reconfigurations are not hitless, we note that the link availability under RADWAN does not suffer significantly as very few links undergo rapid changes in capacity. This is confirmed by Figure 10 which shows the distribution of the number of capacity reconfigurations observed during the simulation period.

As expected, in the absence of catastrophic optical events ($\text{SNR} < 3$) during the simulation period, RADWAN-HITLESS allows links to be available all the time by instantly adapting the link capacity to the lower or higher SNR. We also find that SWAN-150 achieves the same availability as SWAN-100 in our simulation.

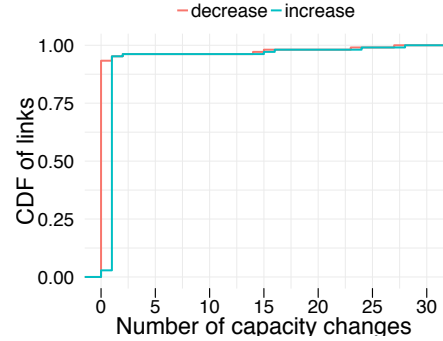


Figure 10: Distribution of the number of capacity reconfigurations occurring per link in the network. We note that only 6% of the links change their capacities more than once in the simulation period.

7 DISCUSSION

In this section, we consider future directions of rate adaptive networks and suggest means of achieving hitless capacity change. We then discuss the impact of underlying fiber length on dynamic capacity links and the cost of operating them.

7.1 Hitless Capacity Change

BVTs and dependency graphs. Dependency graphs [21, 24] are a seminal technique used for consistent network updates [10]. To perform consistent updates, an old and a new network state is specified such that a routing change is performed only when safe to do so. However, to change the capacity of a link e , carrying flow f before and after the capacity reconfiguration, dependency graphs perform poorly since no alternative path is specified for f .

RADWAN manages link flaps by computing an intermediate routing state for flows during reconfiguration.

As such, RADWAN specifies a two-step dependency graph: in order for a scheduled link flap to be activated, the affected traffic is rerouted beforehand.³ Because of the benevolent nature of SNR in our dataset, coupled with the churn minimization of Algorithm 1, RADWAN jointly activates all link flaps. In more volatile SNR scenarios, RADWAN can be set to activate link flaps over multiple dependent iterations. We conjecture that such intermediate consistency methods can eventually be phased out once hitless capacity changes become production ready, as discussed in the next section.

Towards hitless capacity change. BVTs are not yet optimized to handle the latency of a modulation change. State-of-the-art BVTs can only change the link modulation after bringing the module to a lower power state. This translates to a link flap for higher layer protocols. The duration of such link failures is a challenge in the deployment of dynamic capacity links in production networks. To quantify this, we obtain an evaluation board of the Acacia AC400 bandwidth variable transceiver [1]. This is the same module which is integrated in the switch linecard used as part of our testbed in Section 5. Since the evaluation board exposes an API to program the transceiver, we use it to understand the modulation change procedure. We change the link’s modulation 200 times from QPSK to 16QAM and analyze the time taken.

Figure 11a shows the AC400 bandwidth variable transceiver module. We observe that the average downtime of the link undergoing capacity change is 68 seconds, similar to the observation made in Section 5. We investigate the cause of latency in capacity reconfiguration and find that the majority of this time is associated with turning the laser back on after reprogramming the transceiver module. We plot the distribution of time taken to change modulation without turning off the laser (Figure 11b) and find that it only takes approximately 35 ms on average. This suggests an opportunity to strive towards hitless capacity changes in the fiber.

7.2 Cost and Distance

One of the key benefits of deploying bandwidth variable links is their cost savings. While the exact cost of individual transceivers is highly dependent on bulk discounts offered by device manufacturers, conversations with industrial partners suggest that the cost of BVTs is on par with the cost of 100 Gbps static transceivers. Due to comparable costs of the two transceivers, operators are increasingly adopting BVTs even though their modulation format is programmed only a handful of times.

³OWAN [20] also deals with consistent cross-layer reconfiguration in WANs, but it is designed for Reconfigurable Optical Add-Drop Multiplexers, where wavelengths are exclusively *either* activated or deactivated: link flaps due to BVTs are not considered.

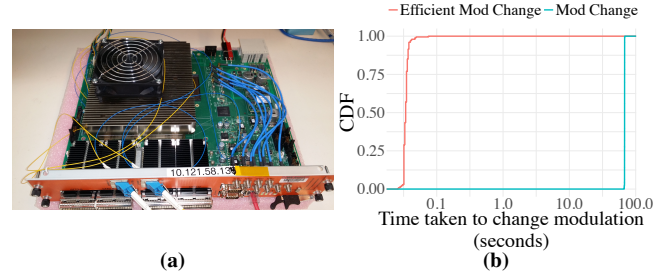


Figure 11: (a) AC400 BVT evaluation board to analyze modulation change latency. (b) CDF of the time taken to change modulation (capacity) of a fiber link using the BVT. Link capacity changes take 68 seconds, on average, but we demonstrate ways to change the modulation efficiently, so that it takes only 35 milliseconds.

RADWAN allows operators to take advantage of BVTs by enabling higher data rates and consequently reducing the dollars per gigabit (\$/Gb) value of traffic in the network. Using the distribution of potential link capacities (Figure 1) enabled by Acacia BVTs and the \$/Gb cost of sending traffic, we estimate that RADWAN provides an overall cost saving of at least 32% over the state-of-the-art.

A caveat of using higher order modulations is that they limit the distance light can travel in fiber. This is because higher number of symbols (as in 8 QAM and 16QAM) in the modulation format reduces the minimum distance between adjacent symbols, making the transmission more prone to distortion as the signal traverses longer distances [9].

As mentioned in §1, QPSK modulation format supports data rates of 100 Gbps for distances upto 3,000 km, 8QAM allows 150 Gbps for distances up to 2,100 km, and 16QAM allows 200 Gbps for distances up to 800 km. We analyzed the fiber distances in our WAN and found that the majority of our fiber paths are less than 800 km (thus capable of supporting 16QAM) and only a small percentage of paths are longer than 2,100 km. While our current proposal did not take fiber length into account, we believe it can be extended to incorporate distance as a constraint.

8 RELATED WORK

Our work builds on several lines of related research as categorized below.

Optical and IP layer orchestration. Singh *et al.* [27] recently analyzed the SNR of links in a large North American backbone over a period of 2.5 years and proposed adapting link capacities to the SNR optical channels. We extend their study period to three years, and, at the same time, broaden their initial measurement and testbed quantifications. We also propose a centralized TE controller system RADWAN and evaluate the interaction between dynamic capacity links and

IP layer flows with simulations at scale and in a realistic testbed. The study by Jin *et al.* [20] on cross-layer optimization between IP and optical layers wavelengths is similar in spirit to our motivation of bridging the gap between optical and IP layers. In their work, Jin *et al.* show the reconfiguration of wavelengths provides latency gains for deadline-driven bulk transfers, also providing a competitive analysis of scheduling single-hop transfers in [19]. But their work keeps the capacity of each wavelength static. In contrast, our work focuses on the reconfiguration of the *capacity* of wavelengths, without the migration of wavelengths across links. In addition, we provide measurements from an operational backbone and argue for changing link capacities with a focus on throughput and reliability. An interesting future direction would be to study the throughput and latency gains of a combination of the two proposals: a fully programmable WAN topology where both capacities and placement of wavelengths on fiber is informed by the centralized TE.

WAN measurements. Govindan *et al.* [14] study 100 failure events across two WANs and data center networks, offering insights into the challenges of maintaining high levels of availability for content providers. Although they do not isolate optical layer failures, they report on root causes of failures, including optical transmitters. We complement their work by focusing on optical layer failures. Ghobadi *et al.* study Q-factor data from Microsoft’s optical backbone [12, 13] and provide insights into the data. Our work complements their analysis on several fronts. First, we make a deeper dive into the impact of temporal changes of SNR on link capacities in terms of capacity gain, availability gains, and realistic throughput gains. Second, we propose and build the system infrastructure required to achieve capacity variable links and benchmark the throughput gains using realistic IP level data. Third, we build a comprehensive testbed and evaluate the impact of capacity reconfiguration, as well as amplifiers, on the path. Our work closes the loop for enabling capacity variable links. Similarly, Filer *et al.* [7] studied the deployed optical infrastructure of Microsoft’s backbone; they discuss the benefits of optical elasticity, express a long-term goal of unifying the optical control plane with routers under a single Software Defined Network controller and recognize YANG [4] and SNMP as potential starting points for a standard data model and control interface between the optical layer and the WAN traffic controller. In this work, we explore how programmability in the optical layer yield throughput gains, and we present a cross-layer WAN traffic controller for dynamic capacity links. Marian *et al.* [25] focused on IP and TCP layer measurements, such as packet loss and packet inter-arrival times, on fiber optics spans. In contrast, we capture failures in the optical layer using failure tickets.

Hardware feasibility studies. Yoshida *et al.* [31, 32] studied the use of 12.5 GHz spectrum slices to allocate bandwidth variable connections to improve the spectrum usage. Although their works did not consider real-time adjustment of the capacity, it provided the foundation for the feasibility of building the necessary hardware with variable bandwidth capabilities—the enabler of our work. We use real-world measurements and build a system that fills the gap between optical and IP layers. Fischer *et al.* [8] and Teipen *et al.* [29] efforts to commercialize higher-speed optical transmission have demonstrated the need for advanced modulation formats, several of which require similar transceiver hardware architecture. Their work showed that adaptive transceivers can be built to support a number of possible operational configurations, but they did not employ a real-time reconfiguration mechanism. In contrast, we discuss the advantages of reconfigurable capacities in real-time based on live SNR measurements.

9 CONCLUSION

In this work, we quantify the throughput and reliability benefits of rate adaptive wide area networks. Our analysis of the SNR of over 8,000 links in an optical backbone for a period of three years shows that the capacity of 64% of the IP links can be increased by 75 Gbps, yielding an overall throughput gain of 134 Tbps. Furthermore, 25% of link failures can be avoided by reducing the transmission rate to 50 Gbps from 100 Gbps. To leverage these benefits, we present RADWAN, a traffic engineering system that dynamically adapts link rates to enhance network throughput and availability. We evaluate RADWAN in a testbed with 1,540 km optical fiber and also simulate throughput and availability gains at scale. By simulating the traffic demand and failures of four random days, we show RADWAN can achieve 40% higher throughput than SWAN. We also address the challenge of the hardware delay in modifying a link’s capacity. We analyze the cause of this delay in current optical transceivers and propose a potential solution to reduce this delay from over a minute to a few milliseconds.

ACKNOWLEDGEMENTS

We would like to thank Victor Bahl, Jamie Gaudette, Jeff Cox, Liban Buni, and Kelly Becker for enabling this study. We thank Mike Pan, Bradford Wright, Urvis Panchal, Megha Sinha, Aditya Bhiday, Rick Ruta, Devin Thorne, and Ratul Mahajan for helpful discussions. We also thank our shepherd David Oran and the anonymous SIGCOMM reviewers for their feedback.

REFERENCES

- [1] Acacia Communications. 2015. Acacia Bandwidth Variable Transceiver Module. <http://ir.acacia-inc.com/phoenix.zhtml?c=254242&p=irol-newsArticle&ID=2103147>. (March 2015).
- [2] Arista Networks. 2017. Python client for Arista eAPI. <https://github.com/arista-eosplus/pyeapi>. (Dec. 2017).
- [3] John C. Bicket. 2005. *Bit-rate Selection in Wireless Networks*. Master's thesis. Massachusetts Institute of Technology.
- [4] Martin Bjorklund. 2010. YANG - A Data Modeling Language for the Network Configuration Protocol (NETCONF). RFC 6020. (Oct. 2010).
- [5] Steven Diamond and Stephen P. Boyd. 2016. CVXPY: A Python-Embedded Modeling Language for Convex Optimization. *Journal of Machine Learning Research* 17 (2016), 83:1–83:5.
- [6] E. Rosen, Y. Rekhter. 2006. BGP/MPLS IP Virtual Private Networks (VPNs). RFC 4364. (Feb. 2006).
- [7] Mark Filer, Jamie Gaudette, Monia Ghobadi, Ratul Mahajan, Tom Issenbuth, Buddy Klinkers, and Jeff Cox. 2016. Elastic Optical Networking in the Microsoft Cloud. *Journal of Optical Communications and Networking* 8, 7 (July 2016), A45–A54.
- [8] J. K. Fischer, S. Alreesh, R. Elschner, F. Frey, M. Nölle, C. Schmidt-Langhorst, and C. Schubert. 2014. Bandwidth-Variable Transceivers based on Four-Dimensional Modulation Formats. *Journal of Lightwave Technology* 32, 16 (Aug 2014), 2886–2895.
- [9] J. K. Fischer, S. Alreesh, R. Elschner, F. Frey, M. Nölle, and C. Schubert. 2013. Bandwidth-variable transceivers based on 4D modulation formats for future flexible networks. In *39th European Conference and Exhibition on Optical Communication (ECOC 2013)*. 1–3. <https://doi.org/10.1049/cp.2013.1366>
- [10] Klaus-Tycho Foerster, Stefan Schmid, and Stefano Vissicchio. 2016. Survey of Consistent Network Updates. *CoRR* abs/1609.02305 (Sept. 2016).
- [11] Saul I. Gass and Arjang A. Assad. 2006. *An Annotated Timeline of Operations Research: An Informal History*. Springer-Verlag New York, Inc., Secaucus, NJ, USA.
- [12] Monia Ghobadi, Jamie Gaudette, Ratul Mahajan, Amar Phanishayee, Buddy Klinkers, and Daniel Kilper. 2016. Evaluation of Elastic Modulation Gains in Microsoft's Optical Backbone in North America. In *Optical Fiber Communication Conference*. Optical Society of America, M2J.2.
- [13] Monia Ghobadi and Ratul Mahajan. 2016. Optical Layer Failures in a Large Backbone. In *Internet Measurement Conference*. ACM.
- [14] Ramesh Govindan, Ina Minei, Mahesh Kallahalla, Bikash Koley, and Amin Vahdat. 2016. Evolve or Die: High-Availability Design Principles Drawn from Googles Network Infrastructure. In *SIGCOMM Conference*. ACM.
- [15] Daniel Halperin, Wenjun Hu, Anmol Sheth, and David Wetherall. 2010. Predictable 802.11 packet delivery from wireless channel measurements. *SIGCOMM Comput. Commun. Rev.* 41, 4 (Aug. 2010), 12.
- [16] Chi-Yao Hong, Srikanth Kandula, Ratul Mahajan, Ming Zhang, Vijay Gill, Mohan Nanduri, and Roger Wattenhofer. 2013. Achieving High Utilization with Software-driven WAN. *SIGCOMM Comput. Commun. Rev.* 43, 4 (Aug. 2013), 15–26.
- [17] Sushant Jain, Alok Kumar, Subhasree Mandal, Joon Ong, Leon Poutievski, Arjun Singh, Subbaiah Venkata, Jim Wanderer, Junlan Zhou, Min Zhu, Jon Zolla, Urs Hölzle, Stephen Stuart, and Amin Vahdat. 2013. B4: Experience with a Globally-deployed Software Defined Wan. *SIGCOMM Comput. Commun. Rev.* 43, 4 (Aug. 2013), 3–14.
- [18] Virajith Jalaparti, Ivan Bliznets, Srikanth Kandula, Brendan Lucier, and Ishai Menache. [n. d.]. Dynamic Pricing and Traffic Engineering for Timely Inter-Datcenter Transfers. In *Proceedings of the 2016 ACM SIGCOMM Conference*.
- [19] Su Jia, Xin Jin, Golnaz Ghasemiesfeh, Jiaxin Ding, and Jie Gao. 2017. Competitive analysis for online scheduling in software-defined optical WAN. In *INFOCOM*. IEEE.
- [20] Xin Jin, Yiran Li, Da Wei, Siming Li, Jie Gao, Lei Xu, Guangzhi Li, Wei Xu, and Jennifer Rexford. 2016. Optimizing Bulk Transfers with Software-Defined Optical WAN. In *SIGCOMM Conference*. ACM.
- [21] Xin Jin, Hongqiang Harry Liu, Rohan Gandhi, Srikanth Kandula, Ratul Mahajan, Ming Zhang, Jennifer Rexford, and Roger Wattenhofer. 2014. Dynamic Scheduling of Network Updates. *SIGCOMM Comput. Commun. Rev.* 44, 4 (Aug. 2014), 539–550.
- [22] Srikanth Kandula, Ishai Menache, Roy Schwartz, and Spandana Raj Babbula. [n. d.]. Calendaring for Wide Area Networks. In *Proceedings of the 2014 ACM Conference on SIGCOMM*.
- [23] Hongqiang Harry Liu, Srikanth Kandula, Ratul Mahajan, Ming Zhang, and David Gelernter. 2014. Traffic Engineering with Forward Fault Correction. *SIGCOMM Comput. Commun. Rev.* 44, 4 (Aug. 2014), 527–538.
- [24] Ratul Mahajan and Roger Wattenhofer. 2013. On consistent updates in software defined networks. In *HotNets*. ACM.
- [25] T. Marian, D.A. Freedman, K. Birman, and H. Weatherspoon. 2010. Empirical characterization of uncongested optical lambda networks and 10GbE commodity endpoints, In DSN. *DSN*.
- [26] Andrew McGregor and Derek Smithies. 2010. Rate Adaptation for 802.11 Wireless Networks: Minstrel. <http://blog.cerowrt.org/papers/minstrel-sigcomm-final.pdf>.
- [27] Rachee Singh, Monia Ghobadi, Klaus-Tycho Foerster, Mark Filer, and Phillipa Gill. 2017. Run, Walk, Crawl: Towards Dynamic Link Capacities. In *HotNets*. ACM.
- [28] Spirent Communications. 2018. Spirent TestCenter. <https://www.spirent.com/Products/TestCenter>. (Jan. 2018).
- [29] Brian Thomas Teipen, Michael Eiselt, Klaus Grobe, and Jörg-Peter Elbers. 2012. Adaptive Data Rates for Flexible Transceivers in Optical Networks. 7 (05 2012).
- [30] Mythili Vutukuru, Hari Balakrishnan, and Kyle Jamieson. 2009. Cross-layer Wireless Bit Rate Adaptation. *SIGCOMM Comput. Commun. Rev.* 39, 4 (Aug. 2009), 3–14.
- [31] Y. Yoshida, A. Maruta, K. i. Kitayama, M. Nishihara, T. Tanaka, T. Takahara, J. C. Rasmussen, N. Yoshikane, T. Tsuritani, I. Morita, S. Yan, Y. Shu, Y. Yan, R. Nejabati, G. Zervas, D. Simeonidou, R. Vilalta, R. Muñoz, R. Casellas, R. Martinez, A. Aguado, V. Lopez, and J. Marhuenda. 2015. SDN-Based Network Orchestration of Variable-Capacity Optical Packet Switching Network Over Programmable Flexi-Grid Elastic Optical Path Network. *Journal of Lightwave Technology* 33, 3 (Feb 2015), 609–617.
- [32] Y. Yoshida, A. Maruta, K. Kitayama, M. Nishihara, T. Tanaka, T. Takahara, J. C. Rasmussen, N. Yoshikane, T. Tsuritani, I. Morita, S. Yan, Y. Shu, M. Channegowda, Y. Yan, B. R. Rofoee, E. Hugues-Salas, G. Saridis, G. Zervas, R. Nejabati, D. Simeonidou, R. Vilalta, R. Muñoz, R. Casellas, R. Martinez, M. Svaluto, J. M. Fabrega, A. Aguado, V. Lopez, J. Marhuenda, O. G. de Dios, and J. P. Fernandez-Palacios. 2014. First international SDN-based network orchestration of variable-capacity OPS over programmable flexi-grid EON. In *OFC 2014*. 1–3.