# ARROW: Restoration-Aware Traffic Engineering

Zhizhen Zhong
Massachusetts Institute of Technology

Manya Ghobadi
Massachusetts Institute of Technology

Alaa Khaddaj
Massachusetts Institute of Technology

Jonathan Leach
Facebook

Yiting Xia
Max Planck Institute for Informatics

Ying Zhang
Facebook

## ABSTRACT

Fiber cut events reduce the capacity of wide-area networks (WANs) by several Tbps. In this paper, we revive the lost capacity by *reconfiguring* the wavelengths from cut fibers into healthy fibers. We highlight two challenges that made prior solutions impractical and propose a system called ARROW to address them. First, our measurements show that contrary to common belief, in most cases, the lost capacity is only *partially restorable*. This poses a *cross-layer* challenge from the Traffic Engineering (TE) perspective that has not been considered before: "*Which IP links should be restored and by how much to best match the TE objective?*" To address this challenge, ARROW's restoration-aware TE system takes a set of partial restoration candidates (that we call LotteryTickets) as input and proactively finds the best restoration plan. Second, prior work has not considered the *reconfiguration latency* of amplifiers. However, in practical settings, amplifiers add tens of minutes of reconfiguration delay. To enable fast and practical restoration, ARROW leverages optical noise loading and bypasses amplifier reconfiguration altogether. We evaluate ARROW using large-scale simulations and a testbed. Our testbed demonstrates ARROW's end-to-end restoration latency is eight seconds. Our large-scale simulations compare ARROW to the state-of-the-art TE schemes and show it can support 2.0×–2.4× more demand without compromising 99.99% availability.

## CCS CONCEPTS

• **Networks** → **Wide area networks**; **Traffic engineering algorithms**; *Network reliability*; *Layering*; Network simulations; Network experimentation; Network measurement; • **Computer systems organization** → **Availability**; • **Mathematics of computing** → *Probabilistic algorithms*;

## KEYWORDS

Wide-area networks, Traffic engineering, Optical restoration, Randomized rounding, Network optimization

## 1 INTRODUCTION

Fiber cuts are undesirable events in Wide-Area Networks (WANs) because (*i*) each fiber carries several Tbps of traffic, and (*ii*) fiber cuts tend to take a long time to repair. Our analysis of failure tickets at Facebook shows fiber cuts account for 67% of total downtime and 50% of the fiber cut events take over nine hours to repair (§2).

Today's service providers cope with the loss of capacity caused by fiber cut events by over-provisioning the network. In particular, to protect from massive packet loss, WAN operators *pre-allocate* extra capacity for failover paths using (*i*) failure-aware Traffic Engineering (TE) [17, 40, 48, 63, 79] and (*ii*) optical path protection [14, 19, 59, 68, 80, 81, 84]. In such techniques, when fiber cuts occur, traffic is automatically shifted from failed IP links to pre-allocated backup paths. We argue that pre-allocating paths is unnecessarily expensive because when a fiber is cut, the router ports and transponders associated with that fiber are still usable.

A more attractive solution is to *reconfigure* the cut fiber's wavelengths to healthy fibers enabling the transponders and router ports associated with the cut fiber to carry traffic while the fiber itself is out of commission. This idea is called *optical restoration* and was proposed two decades ago [30]. Despite several follow up papers [32, 45, 49, 52, 56, 72, 74, 82] and the presence of commercially available devices capable of wavelength reconfiguration, such as Reconfigurable Optical Add Drop Multiplexers (ROADMs) [2, 4], the deployment of prior proposals of optical restoration at scale has several challenges, as we describe next.

Our measurements of a global WAN with over 200,000 IP links and 1,000 optical fiber links show that in practice, 62% of fiber cuts have to be *partially* restored because the remaining fibers do not have enough available spectrum to host all the wavelengths of the cut fiber. Hence, a practical restoration system must be able to choose which IP links to restore partially and by how much.

We demonstrate that when full restoration is not possible, simply maximizing the restored bandwidth from the optical layer's perspective, without considering IP links' traffic demand, leads to sub-optimal throughput. Hence, an IP/optical *cross-layer TE* is needed to ensure the partially restored capacity is carefully allocated across IP links and is efficiently utilized to match the current traffic demand.

However, today's TE schemes [17, 40, 42, 47, 48, 58, 60, 63, 77, 79] do not consider optical restoration for fiber cuts. Instead, they consider fiber cuts as fatal events: an IP link is down until the fiber is repaired. To incorporate partial restoration into the TE, we propose a *restoration-aware* TE system, called Agile RestoRation of Optical Wavelengths (ARROW). ARROW grapples with the *algorithmic* challenges of formulating a restoration-aware TE, such as how to find

the best partial restoration plan while optimizing network throughput, and also with *system-level* challenges, such as how to reduce the end-to-end restoration latency.

In particular, Arrow's TE formulation accounts for partial restoration candidates for IP links during hypothetical fiber cut scenarios and plans according to the best restoration plan *proactively*. We show when full restoration is not possible, many partial restoration candidates can maximize the total restored bandwidth from the optical layer's perspective, but not all of them maximize network throughput from the TE perspective. Taking all of them into account in a joint IP/optical optimization is computationally infeasible, while taking only one of them is sub-optimal.

Arrow solves this problem in a two-stage approach. The first stage involves an *offline* analysis of the available fiber spectrum to find a set of potential restoration candidates, which we call LotteryTickets (§3.2). Each LotteryTicket represents one possible restoration candidate without taking instantaneous traffic demand into account. The LotteryTickets serve as an abstraction between the optical and IP layers. The second stage solves an *online* TE formulation for the current traffic demand to find the *winning* LotteryTicket for each hypothetical fiber cut scenario (§3.3). The TE formulation finds the appropriate tunnel allocations and restoration plans *proactively*, before fiber cuts happen, enabling the network to react quickly when a particular fiber cut happens.

An important practical consideration is the end-to-end reconfiguration latency. Although ROADMs [4] have been ubiquitously deployed in our WAN, reconfiguring a set of wavelengths from a cut fiber to healthy ones results in optical power instability on the *amplifiers* of the new optical path(s). Such power excursions will cause packet loss until all the amplifiers adjust their optical gain, a process that takes several minutes in practice [16, 67, 90, 92]. To address this challenge, Arrow leverages a recently commoditized device called the Amplified Spontaneous Emission (ASE) noise source [83] to bypass the amplifier reconfiguration time altogether, reducing the end-to-end restoration latency from tens of minutes to eight seconds. Prior work showed this device is currently being deployed in large-scale WANs to facilitate wavelength installation [33, 35, 54].

To evaluate Arrow, we build a production-level testbed with 4 ROADM sites, 34 amplifiers, and over 2000 km fiber that faithfully emulate part of our production backbone. Using this testbed, we demonstrate the feasibility of reconfiguring 2.8 Tbps IP capacity (14 wavelengths) within eight seconds. To evaluate the impact of Arrow on throughput and availability, we conduct extensive simulations, comparing state-of-the-art TE algorithms (TeaVaR [17], FFC [63], and ECMP [21]) with Arrow. Our simulations show Arrow supports between 2.0×–2.4× more demand without compromising 99.99% availability. Moreover, we demonstrate that Arrow sustains the same throughput at the same availability level while requiring 2.8× fewer router ports and optical transponders.

## 2 BACKGROUND AND MEASUREMENTS

To motivate our work, we investigate the impact of fiber cuts in a subset of Facebook's WAN with more than 200,000 IP links and 10,000 optical wavelengths traversing 1,000 optical fiber cables across the world.
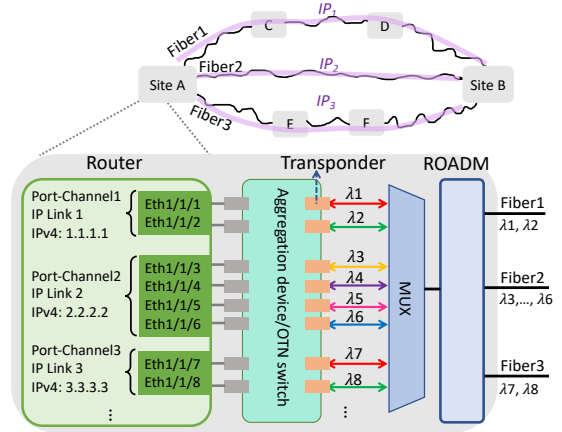


Figure 1: Mapping between IP links and wavelengths.

### 2.1 Overview

Figure 1 illustrates how IP links and optical wavelengths are mapped onto each other in Facebook. As shown, several router ports are grouped into one *port-channel*. Each port-channel represents an IP link and carries several Tbps of traffic via *multiple wavelengths*. Flows destined for an IP address are load-balanced across all the interfaces of a port-channel [9, 25, 46]. The aggregation device aggregates multiple grey router ports into tunable Dense Wavelength-Division Multiplexing (DWDM) transponders.[1] For simplicity of representation, the figure shows a 1-to-1 mapping between router ports and transponders, but a real deployment is more complex. The mapping between wavelengths and fibers is configured in the ROADM. ROADMs can dynamically reconfigure wavelengths to map to any fiber [1, 3, 37], but as we show in this paper, using this feature is not without challenges.

Today, once a fiber is cut, router ports and transponders associated with the failed fiber become unusable and sit idle until the fiber is repaired. However, several *healthy* fibers are often available to reconfigure the wavelengths traversing the cut fiber by reconfiguring the ROADMs on the fiber path. For example, in Fig. 1, if fiber 2 is cut, IP link 2 goes down, causing wavelengths λ3 to λ6, their corresponding transponders, and their router ports (Eth1/1/3 to Eth1/1/6) to become idle. Our goal is to reconfigure these four idle wavelengths (or some of them) on fiber 1 and/or fiber 3. We refer to fiber 1 and fiber 3 as *surrogate fibers*. The decision of which surrogate fiber to choose for each port-channel (i.e., IP link) and how much capacity to restore depends on several factors discussed later in the paper.

Figure 2 presents a high-level example of Arrow in action. The top row shows the network in a healthy state. The first and last columns show the optical and IP-layer views of the network, respectively. The middle column represents the mapping between the two layers. Note the purple IP link between A and C in Figs. 2(b) and (c): even though there are no *direct fibers* between A and C, the provider configured site D to pass through the light between A and
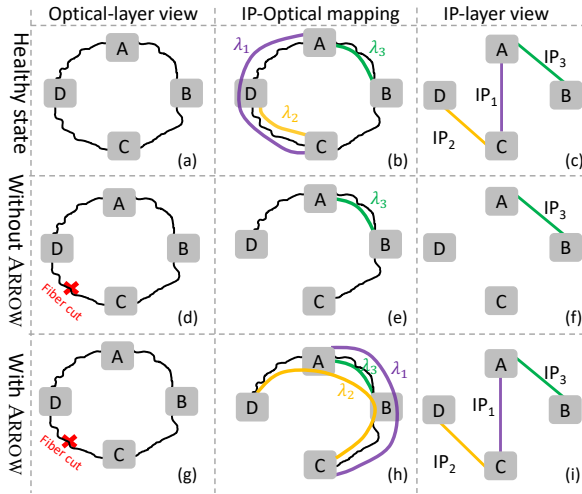
Figure 2: ARROW restores the IP-layer view by reconfiguring wavelengths traversing the cut fiber (i.e., $\lambda_1$ and $\lambda_2$) into healthy ones.



Figure 3: Analysis of 600 failure tickets.



Figure 4: Impact of fiber cuts on IP layer capacity.

C entirely in the optical domain. As a result, the light is not terminated on intermediate hop D, and from the IP layer's perspective, $IP_1$ is a *direct IP link* between nodes A and C.

The second row in Fig. 2 shows the status of the network after a fiber cut without ARROW. As shown in Figs. 2(e) and (f), $IP_1$ and $IP_2$ become unavailable because they were traversing the cut fiber. Consequently, the IP layer has to operate with reduced capacity while the router ports corresponding to $IP_1$ and $IP_2$ sit idle. In contrast, ARROW (third row) restores the IP-layer view back to the healthy state by reconfiguring the wavelengths corresponding to $IP_1$ and $IP_2$ to traverse healthy fibers (through node B), as shown in Figs. 2(h) and (i).

## 2.2 Impact of Fiber Cuts on IP Capacity

We begin by studying 600 WAN-related failure tickets in Facebook over a period of three years (March 2016–June 2019). For each ticket, we record the duration of the failure and its root cause.

Fig. 3(a) plots the CDF of the mean time to repair for all tickets categorized by their root cause. It shows that 50% of the fiber cut events last longer than nine hours, and 10% last over a day. Fig. 3(b) shows the percentage of downtime for each category. As shown, the duration of fiber cut events accounts for 67% of the total downtime. Note that a fiber cut can occur for a variety of reasons including accidental damage by construction workers, aerial poles falling, extreme weather conditions, and fiber being chewed by animals.

To quantify how much IP-layer capacity is lost because of fiber cuts, we dig deeper into the fiber-related failure tickets. We find that, on average, 16 fiber cut events happen every month. Given that Facebook's datacenter sites have multiple fibers between them (see Fig. 1) a fiber cut will take away some capacity between site-pairs. We study the impact of fiber cuts on all site-pairs in Facebook. Fig. 4(a) shows the time series of lost capacity between four site-pairs that suffered the most capacity loss between 2017 and 2018. Each peak in the figure represents a fiber cut, resulting in several Tbps of capacity loss between site pairs.
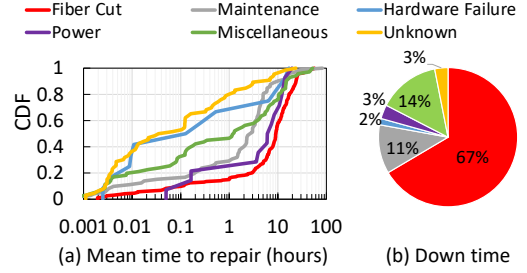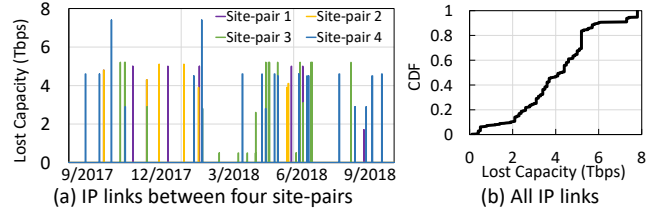
Fig. 4(b) shows the CDF of lost capacity on all IP links caused by fiber cuts during the entire three years of our measurement. We observe that each fiber cut event resulted in the loss of up to 8 Tbps of IP capacity. This massive loss of capacity motivates us to investigate the potential of wavelength reconfiguration for restoring failed IP links.

## 2.3 Partial Restoration

An immediate question to answer is: "Is there enough room in the optical domain such that for every fiber cut scenario, we can reconfigure all affected wavelengths to healthy fibers?" The answer depends on the number of wavelengths that are already provisioned on fibers. Fig. 5(a) shows the CDF of the spectrum utilization (number of provisioned wavelengths divided by the total number of available wavelengths) of Facebook's fibers. The figure shows that 95% of fibers have a spectrum utilization less than 60%. This means 95% of the fibers have 40% spare room for ARROW's wavelength reconfiguration.

Note that the *usable spectrum* for wavelength reconfiguration is usually smaller than the *available spectrum* of each fiber link. This is because, in optical domain, a wavelength's frequency must remain the same throughout the entire fiber path. This property is called *wavelength continuity constraint* in optical networking literature [12, 20, 64, 88]. For example, as shown in Fig. 5(b), although the three fiber links (fiber DA, fiber AB, and fiber BC) all have 75% of their spectrum *available* (spectrum utilization is 25%), it turns out that only 25% of the spectrum is *usable* for reconfiguring $\lambda_4$. If the failed IP link between nodes C and D contains more than one wavelength, it will result in partial restoration, whereby only one wavelength ($\lambda_4$) can be restored.

To quantify the amount of partial restoration in Facebook, we define the *restoration ratio* of a fiber $\phi$ as $U_\phi = \frac{W'_\phi}{W_\phi}$, where $W_\phi$ is the provisioned bandwidth capacity (in Gbps) in healthy state and $W'_\phi$ is the restorable bandwidth capacity after $\phi$ is cut. Each fiber in
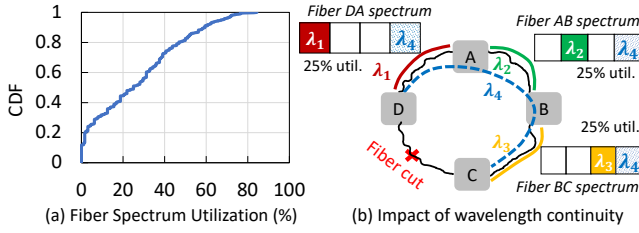
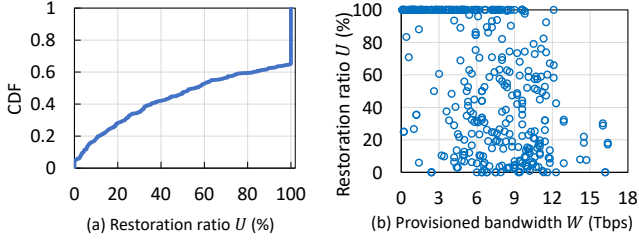**Figure 5: Spectrum utilization of Facebook's fibers.**



**Figure 6: Restoration ratio of Facebook's fibers.**

Facebook carries several wavelengths, $\lambda \in \Lambda$. Hence, $W_\phi = \Sigma_\lambda \beta_\lambda$, where $\beta_\lambda$ is the bandwidth capacity of wavelength $\lambda$. $W'_\phi = \Sigma_\lambda \beta'_\lambda$, where $\beta'_\lambda$ is the *restorable* bandwidth capacity of wavelength $\lambda$. To calculate $W'_\phi$ we iterate over every wavelength $\lambda$ on fiber $\phi$ and check whether we can reconfigure it using the same frequency and modulation on any of the fibers adjacent to $\phi$. Most of the time, the modulation of all wavelengths can be kept the same (details in Appendix A.1). But the frequency of some wavelengths must be tuned to avoid frequency collisions with existing wavelengths already working on surrogate fiber paths using tunable transponders which are already deployed in Facebook. Even with frequency tuning, in some cases, we may not be able to find a common available wavelength frequency end-to-end on all of the fibers along the new fiber path, as different fibers may not have an overlapping available spectrum due to the wavelength continuity constraint. When there is no available frequency, the wavelength is not reconfigurable, and its restorable bandwidth becomes zero, resulting in $W'_\phi < W_\phi$. Hence fiber $\phi$ becomes *partially restorable*.

For instance, in Fig. 2(b), fiber DC has two wavelengths ($\lambda1$ and $\lambda2$) in its healthy state. Fig. 2(h) shows the restoration path traverses three surrogate fibers (DA, AB, BC). Hence, both $\lambda1$ and $\lambda2$ frequencies must be available on all three fibers. If one of the fibers is already occupying the same frequency as $\lambda2$, ARROW searches for another frequency available across all three fibers. If no such frequency exists, ARROW will not reconfigure $\lambda2$, and the restoration ratio for fiber DC becomes 50%.

To quantify the degree of partial restoration in Facebook, we simulate all single fiber cut scenarios in our WAN and calculate the restoration ratio, $U_\phi$, of each fiber $\phi$. Fig. 6(a) plots the CDF of the restoration ratio (percentage) for all fibers. It shows 34% of fibers are fully restorable, 4% are not restorable at all, and 62% are partially restorable. Further, Fig. 6(b) plots the relationship between the restoration ratio and provisioned bandwidth capacity for all fibers. For fibers with provisioned bandwidth capacity larger than 10 Tbps, the restoration ratio is almost never 100%. In disaster scenarios where multiple fiber cuts can happen, partial restoration
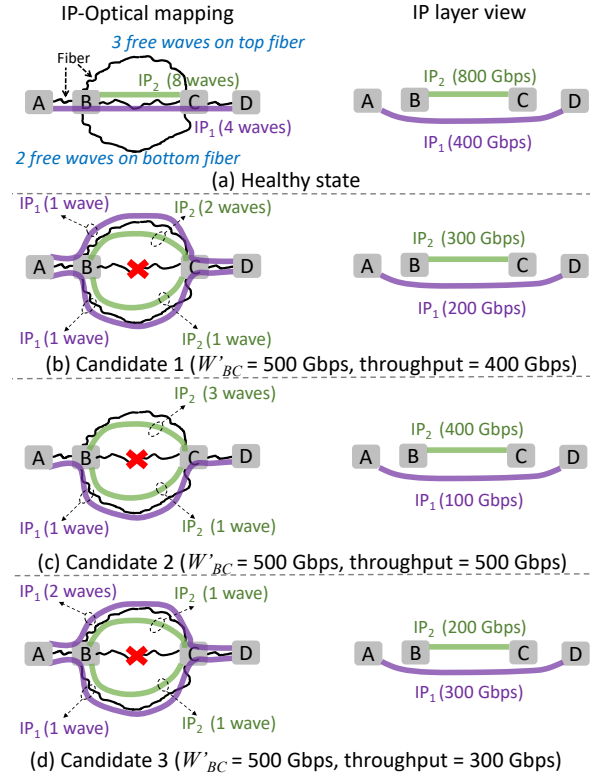


**Figure 7: Illustration of several restoration candidates.**

plays a bigger role because more capacity is lost. Thus, selecting which IP links to recover becomes even more important.

When full restoration is not possible, deciding which IP link to restore and by how much is challenging. Consider the network in Fig. 7(a), with $IP_1$ and $IP_2$ traversing fiber BC. In healthy state, $IP_1$ and $IP_2$ have four and eight wavelengths, respectively. Assuming each wavelength's bandwidth is 100 Gbps, $W_{BC}$ is 1200 Gbps. When fiber BC is cut, its wavelengths can be reconfigured into the top or bottom (or both) surrogate fibers. But suppose the available spectrum is such that the top fiber has three free wavelengths, and the bottom one has two. This means we can only reconfigure five wavelengths, each with 100 Gbps bandwidth ($W'_{BC}$=500 Gbps). However, there are several different options to distribute these wavelengths to each IP link. Fig. 7(b) illustrates restoration candidate 1 where $IP_1$'s restored bandwidth is 200 Gbps (one wavelength restored on the top fiber and another on the bottom fiber). This leaves three wavelengths (two on the top fiber and one on the bottom) for $IP_2$. Two other candidates are illustrated in Figs. 7(c) and (d). Note that all candidate options ultimately contain the *maximum* amount of restorable bandwidth ($W'_{BC}$=500 Gbps). The difference is how they distribute the wavelengths across IP links. From the optical layer's perspective, all of the above restoration candidates are equal. Now, consider the case where $IP_1$'s and $IP_2$'s traffic demands are 100 Gbps and 400 Gbps, respectively. The throughput of restoration candidates 1, 2, and 3 (Figs. 7(b), (c), and (d)) is 400 Gbps, 500 Gbps, and 300 Gbps, respectively. Hence, considering the traffic demand, candidate 2 in Fig. 7(c) is optimal, and the other candidates are sub-optimal.

# 3 HANDLING PARTIAL RESTORATION

Thus far we have shown that, in practice, most optical fibers are not fully restorable. An important question is: *when full restoration is not possible, which partial restoration candidate leads to the best network throughput?* In response, this section explains Arrow's cross-layer approach to handle partial restoration.

## 3.1 High-level Design

We start with a high-level explanation of our design.

**Joint IP/optical formulation is not scalable.** A strawman approach to take partial restorations into account is to jointly optimize the IP-layer's *instantaneous* traffic demand with optical-layer wavelength assignment by formulating the problem as an Integer Linear Program (ILP). However, solving the ILP is not scalable, even for networks of moderate size ($\approx$20 nodes) [40, 91]. In Appendix A.4, we present the joint cross-layer IP/optical optimal formulation. We then report the number of variables and constraints for three network topologies (Table 8) and find that the size of the optimization problem grows massively. Hence, the joint optimization cannot be solved even within several days, making it not suitable for modern TE systems that re-optimize traffic allocation periodically (every few minutes).

**Optical restoration during capacity planning is not optimal.** Instead of taking the instantaneous traffic demand, a common technique in prior work is to plan for optical restoration during the capacity planning phase [13, 23, 30, 31, 39, 40, 79]. While this approach is useful, our evaluations show it is less effective compared to Arrow (§6) because the instantaneous traffic demand can be different from the traffic matrix considered during capacity planning. Moreover, prior proposals still take several minutes, if not hours, to execute as capacity planning is an infrequent event (e.g., months) and does not have a tight execution deadline.

**Centralized TE has a tight execution time.** Modern WANs use centralized TE formulations to reduce congestion and increase efficiency [42, 47, 58, 61]. The TE controller periodically (e.g., every five minutes [6, 42]) adjusts the traffic allocation on network paths to satisfy *the current traffic demand*. As a result, the TE formulation has a tight execution deadline. Our production operators in Facebook allow 5 minutes for TE execution deadline. Hence, to be deployable, Arrow's restoration-aware TE must maintain this tight deadline.

**Abstracting optical layer for TE (§3.2).** To embed *restoration awareness* into centralized TE formulations, we design a two-stage solution. The first stage is executed offline using IP/optical mapping in the WAN, without considering instantaneous traffic demand. As the IP/optical mapping does not change frequently, this stage does not have to be executed as frequently as the TE. The output of this stage is an *abstraction* between the optical layer and the TE. Our abstraction contains a set of partial *restoration candidates*, which we call LotteryTickets.

**Restoration-aware TE formulation (§3.3).** The second stage of our solution is an online restoration-aware TE formulation using LotteryTickets. Using the current traffic demand, our TE formulation computes the best traffic flow allocation and restoration plan
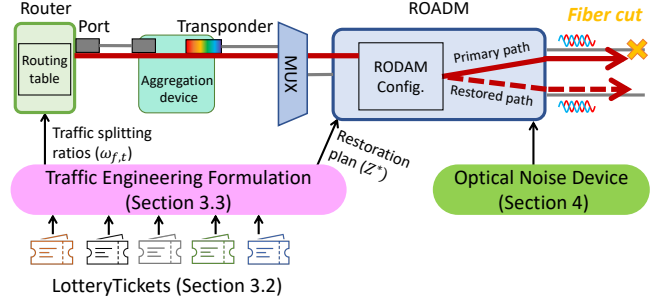


**Figure 8: Arrow's high-level design. Traffic splitting ratios ($\omega_{f,t}$) and restoration plans ($Z^*$) are defined in §3.3.**

for each failure scenario *proactively*, before fiber cuts happen. Finally, Arrow installs the wavelength reconfiguration plans into the ROADMs to be able to react quickly when a particular fiber cut happens. Fig. 8 illustrates Arrow's high-level system design. The LotteryTickets abstraction is explained in Section 3.2. Details of the traffic engineering formulation, including traffic splitting ratios ($\omega_{f,t}$) and restoration plans ($Z^*$), are in Section 3.3.

## 3.2 The LotteryTicket Abstraction

**LotteryTicket design.** The intuition for our LotteryTicket abstraction is shown in Fig. 7: all three restoration candidates are optimal from the optical layer's perspective (maximizing total restorable capacity) but only one, candidate 2, maximizes the *throughput*. In other words, all three candidates are *lottery tickets* with good chances of winning, but depending on the traffic demand, only one is the winner. Given that the traffic matrix changes periodically, simply hard coding one of them into the TE is sub-optimal. Instead, Arrow finds a set $Z$ containing several LotteryTickets and feeds them into the TE as input parameters. The number of LotteryTickets is a configurable parameter. Section 6 evaluates the impact of the number of LotteryTickets on Arrow's performance.

More formally, assume a fiber cut scenario $q$ causes the IP links between $n$ source-destination ROADM pairs to fail. We generate a set $Z$ containing $|Z|$ LotteryTickets, where each ticket $R^{z,q} = \{r_1^{z,q}, ..., r_n^{z,q}\}$ corresponds to the restorable bandwidth capacity of failed IP links indexed from 1 to $n$, where $z \in Z$. For instance, the LotteryTickets corresponding to Fig. 7 (b), (c), and (d) are Ticket$_1$: (200 Gbps, 300 Gbps), Ticket$_2$: (100 Gbps, 400 Gbps), and Ticket$_3$: (300 Gbps, 200 Gbps). Note that all tickets restore a total of 500 Gbps of capacity but the restoration capacity of individual IP links differs across tickets.

**Arrow's LotteryTicket algorithm.** Algorithm 1 shows the pseudocode of our LotteryTicket generation procedure. The algorithm consists of two parts. First, it starts by solving the Routing and Wavelength Assignment (RWA) problem [88], where the ILP is relaxed to LP (line 2). Second, it uses the RWA's solution as the seed for randomized rounding to generate LotteryTickets (lines 3–13).

**Routing and Wavelength Assignment (RWA).** Our **RWA_LP** module takes three inputs: (1) optical-layer network graph $G(\Psi, \Phi)$ with ROADM set $\Psi$ and fiber set $\Phi$; (2) $\overline{\Psi}$ representing the source-destination ROADMs of failed IP links; (3) $\overline{\Lambda}$ representing the wavelength information of failed IP links. The RWA module explores the

**Algorithm 1** Arrow's LotteryTicket Pseudocode

---

1: **procedure** Create LotteryTickets for a fiber cut scenario $q$ that has caused $n$ IP links to be down

   ▷ **Input:** $G(\Psi, \Phi)$: optical-layer network graph with ROADM set $\Psi$ and fiber set $\Phi$

   ▷ **Input:** $\overline{\Psi} = \{< \psi_1^{src}, \psi_1^{dst} >, ..., < \psi_n^{src}, \psi_n^{dst} >\}$: vector representing source-destination ROADM pairs of $n$ failed IP links

   ▷ **Input:** $\overline{\Lambda} = \{\overline{\lambda_1}, ..., \overline{\lambda_n}\}$: vector representing wavelength information (frequency, modulation) of failed IP links

   ▷ **Input:** $\delta$: rounding stride

   ▷ **Input:** $|Z|$: number of LotteryTickets to be generated

   ▷ **Output:** $|Z|$ LotteryTickets $\{R^{z,q}\}$, where $R^{z,q} = \{r_1^{z,q}, ..., r_n^{z,q}\}$ corresponds to the restorable bandwidth capacity of failed IP links, $1 \leq z \leq |Z|$

      ▷ *Find restorable wavelengths* $\Lambda$ *using RWA*

2:    $\Lambda = $ **RWA_LP** $(G(\Psi, \Phi), \overline{\Psi}, \overline{\Lambda})$;   ▷ (Appendix A.2)

3:    **for all** $z \in Z$ **do**

4:        **for all** $1 \leq e \leq n$ **do**

5:            $\lambda_e = \Lambda[e].num\_wavelengths$

              ▷ *Step 1: decide the rounding stride*

6:            $x_1 = $ **randInt**$(1, \delta)$

              ▷ *Step 2: decide rounding up or down*

7:            $x_2 = $ **randFloat**$(0, 1)$

8:            **if** $x_2 < \lambda_e - \lfloor \lambda_e \rfloor$ **then**

                ▷ *round up*

9:                $r_e^{z,q} = min(\lceil \lambda_e \rceil + x_1, \overline{\Lambda[e]}.num\_wavelengths)$

10:          **else**

                ▷ *round down*

11:                $r_e^{z,q} = max(\lfloor \lambda_e \rfloor - x_1, 0)$

           ▷ *Find the restorable capacity according to modulations*

12:         $r_e^{z,q} = r_e^{z,q} \times \Lambda[e].modulation$;

           ▷ *Append the result into LotteryTickets*

13:         $R^{z,q}$.**append** $(r_e^{z,q})$

    **return** $|Z|$ LotteryTickets $\{R^{z,q}\}, z \in Z$

---

| | | |
|---|---|---|
| | $G(V, E)$ | IP-layer network graph with datacenter sites set $V$ and IP links set $E$. |
| | $F = \{f\}$ | Flows aggregated by ingress-egress sites. |
| | $d_f$ | Bandwidth demand of flow $f$. |
| Standard | $c_e$ | Bandwidth capacity of IP link $e \in E$. |
| TE Input | $T_f$ | Set of tunnels for flow $f$, $T_f \subset T$. |
| | $L[t, e]$ | 1 if tunnel $t$ uses IP link $e$ and 0 otherwise. |
| | $Q = \{q\}$ | Considered failure scenarios. |
| | $T_f^q$ | Residual tunnels for flow $f$ under scenario $q$. |

**Table 1: Standard TE input parameters.**

rounding up [78] (lines 7–11).[2] During the rounding process, we make sure the rounded integer is never smaller than zero or goes beyond the initial number of wavelengths for that IP link (lines 9 and 11). Finally, we calculate the amount of restorable bandwidth capacity by multiplying the number of restorable wavelengths by the modulation format (line 12).[3] The algorithm returns $|Z|$ LotteryTickets.

**Handling LotteryTickets' feasibility.** Before Arrow feeds the LotteryTickets into the TE, it performs an additional check to make sure all LotteryTickets are feasible in the optical domain. We add a *feasibility check* module to drop infeasible LotteryTickets that do not meet all the constraints of our RWA formulation (Appendix A.2). More specifically, since the LotteryTickets are generated using a randomized rounding process agnostic to the optical topology, some of them may violate the RWA constraints. Hence, it is necessary to check the feasibility of the generated tickets and filter out the infeasible ones.

### 3.3 Restoration-Aware Traffic Engineering

This section describes Arrow's restoration-aware TE formulation. For clarity of presentation, we use FFC's [63] notation. Note that other techniques that improve over FFC, such as TeaVaR [17] and PCF [48], can also be applied on top of our formulation.

**Standard TE input.** We begin by considering the standard TE input parameters listed in Table 1. Similar to prior work, we model the WAN as a directed graph $G = (V, E)$, where the vertex set $V$ represents the datacenter sites, and edge set $E$ represents the IP links between them. In each time epoch, there is a set of source-destination pairs (or "flows"), where each such pair $f$ is associated with a demand $d_f$, and a fixed set of paths (or "tunnels") $T_f \subset T$ on which its traffic should be routed. Link capacities are given by $C = (c_1, ..., c_{|E|})$ (e.g., in bps). Similar to FFC, Arrow assumes the tunnels are part of the input, but the formulation can be extended to approaches such as $k$-shortest paths, traffic oblivious tunnels [58], or logical sequences [48]. Failure scenarios are denoted by $q \in Q$, including fiber cuts, switch failures, and control plane failures. In this paper, we only consider IP link failures caused by fiber cuts, since optical restoration does not apply to switch or control plane failures. To avoid having an exponential number of failure-related

optical topology of each failed IP link and tries to find a possible surrogate fiber path for each lost wavelength. To avoid the extra latency associated with frequency tuning and/or modulation change, Arrow tries to keep the same modulation and frequency, if at all possible. Otherwise, it finds the best alternative modulation and frequency assignment. Since the RWA is a well-studied problem in the optical networking literature [11, 12, 88, 91], we omit the formulation here and refer the reader to Appendix A.2 for details.

**Generating LotteryTickets with randomized rounding.** The solution of **RWA_LP** is a set $\Lambda$ containing the frequency and modulation of the restorable wavelengths for IP links (line 2). However, because of the ILP to LP relaxation, the number of restorable wavelengths is not always an integer. It turns out this is a blessing in disguise for our LotteryTicket abstraction. We take advantage of the situation by *repeating* a randomized rounding technique [70] to generate $|Z|$ LotteryTickets from the floating point solution (lines 3–13).

To construct each LotteryTicket $R^{z,q}$, we start with the optimal floating-point solution $\Lambda$ (line 5). The rounding process has two probabilistic steps: 1) it decides the rounding stride based on a random integer within $[1, \delta]$ where $\delta$ is an input parameter (line 6); and 2) it decides the rounding direction (up or down) by taking the fractional part of the floating point solution as the probability of

---

| | Table 1 | Standard TE input parameters. |
| --- | --- | --- |
| | $Z^q = \{z\}$ | Set of LotteryTicket indexes under scenario $q$. |
| ARROW Phase I Input Parameters | $r_e^{z,q}$ | Restorable bandwidth capacity for link $e$ under scenario $q$ and LotteryTicket $z$. |
| | $Y_f^{z,q}$ | Restorable tunnels for flow $f$ under scenario $q$ and LotteryTicket $z$. |
| | $M^{z,q}$ | A parameter to bound LotteryTicket $z$'s slack variables under scenario $q$. |
| ARROW Phase I Output | $\Delta_e^{z,q}$ | Slack variable for each edge $e$'s restorable bandwidth capacity $r_e^{z,q}$ under scenario $q$ and LotteryTicket $z$. |

**Maximize:** $\sum_{f \in F} b_f$

**Subject to:**

$$\forall f: \quad \sum_{t \in T_f} a_{f,t} \geq b_f \quad (1)$$
$$\forall e: \quad \sum_{f \in F} \sum_{t \in T_f} a_{f,t} \times L[t,e] \leq c_e \quad (2)$$
$$\forall f: \quad 0 \leq b_f \leq d_f \quad (3)$$
$$\forall f,q,z: \quad \sum_{t \in Y_f^{z,q}} a_{f,t} + \sum_{t \in T_f^q} a_{f,t} \geq b_f \quad (4)$$
$$\forall e,q,z: \quad \sum_{f \in F} \sum_{t \in Y_f^{z,q}} a_{f,t} \times L[t,e] \leq r_e^{z,q} + \Delta_e^{z,q} \quad (5)$$
$$\forall q,z: \quad \sum_{e \in E} \Delta_e^{z,q} \leq M^{z,q} \quad (6)$$

**Table 2: ARROW TE Phase I formulation.**

constraints, we use TeaVaR's probabilistic approach [17] and only consider highly-probable failure scenarios (see §6 for details).

**ARROW's two phase formulation.** Given a set of LotteryTickets, our goal is to find the winning LotteryTicket for each failure scenario and instantaneous traffic matrix. However, this formulation would require solving an ILP (as shown in Table 9 in Appendix A.5) which violates the tight runtime deadline for ARROW TE. To address this challenge, we introduce a two-phase formulation that separates the LotteryTickets selection from the traffic allocation process while keeping both of them in LP form. Phase I selects the winning LotteryTicket for each failure scenario based on the input traffic demand and Phase II uses the winning LotteryTicket to find the best traffic allocation on tunnels. Tables 2 and 3 present ARROW's Phase I and II formulations, respectively.

**Phase I input parameters.** In addition to the standard TE inputs, Phase I formulation takes the following input parameters: a series of LotteryTickets (from Algorithm 1) where the amount of restorable bandwidth capacity for IP link $e$ is given as $r_e^{z,q}$. Moreover, Phase I's input parameters include a set of *restorable tunnels* for flow $f$ under scenario $q$ and LotteryTicket $z$, denoted by $Y_f^{z,q}$. A tunnel is considered restorable if some (or all) of its IP links are restorable during failure scenario $q$. The set $Y_f^{z,q}$ is calculated based on $r_e^{z,q}$. In scenario $q$, if every failed link $e$ that tunnel $t$ traverses is available after restoration (i.e., $\prod_{e \in E} L[t,e] \times r_e^{z,q} > 0$), this tunnel is restorable under scenario $q$ (i.e., $t \in Y_f^{z,q}$). Finally, we introduce a new input parameter called $M^{z,q}$ calculated as $\alpha \times \sum_{e \in E} r_e^{z,q}$ to capture the $\alpha$-fraction of total restorable bandwidth capacity for scenario $q$ and LotteryTicket $z$ (details below).[4]

**Phase I optimization goal and constraints.** We use the same optimization goal as FFC to maximize the network throughput. Constraints (1-3) are standard TE constraints to ensure the following: the sum of the bandwidth of all tunnels of flow $f$ should be larger

---

| | Table 1 | Standard TE input parameters. |
| --- | --- | --- |
| ARROW Phase II Input Parameters | $r_e^{*,q}$ | Winning LotteryTicket's restorable bandwidth capacity for link $e$ under scenario $q$. |
| | $Y_f^{*,q}$ | Winning LotteryTicket's restorable tunnels for flow $f$ under scenario $q$. |
| ARROW Phase II Output | $b_f$ | Total allocated bandwidth for flow $f$. |
| | $a_{f,t}$ | For flow $f$, the allocated bandwidth on tunnel $t \in T_f$. |

**Maximize:** $\sum_{f \in F} b_f$

**Subject to:**

$$\forall f: \quad \sum_{t \in T_f} a_{f,t} \geq b_f \quad (7)$$
$$\forall e: \quad \sum_{f \in F} \sum_{t \in T_f} a_{f,t} \times L[t,e] \leq c_e \quad (8)$$
$$\forall f: \quad 0 \leq b_f \leq d_f \quad (9)$$
$$\forall f,q: \quad \sum_{t \in Y_f^{*,q}} a_{f,t} + \sum_{t \in T_f^q} a_{f,t} \geq b_f \quad (10)$$
$$\forall e,q: \quad \sum_{f \in F} \sum_{t \in Y_f^{*,q}} a_{f,t} \times L[t,e] \leq r_e^{*,q} \quad (11)$$

**Table 3: ARROW TE Phase II formulation.**

than $f$'s allocated bandwidth $b_f$ with Constraint (1); the sum of the bandwidth of all tunnels on a given IP link $e$ should be no larger than the link capacity $c_e$ with Constraint (2); and the allocated bandwidth of flow $f$ should be less than the demand of $f$ with Constraint (3). Constraint (4) considers ARROW's restorable tunnels by ensuring the sum of the bandwidth of both residual *and restorable* tunnels for flow $f$ under failure scenario $q$ is larger than $f$'s allocated bandwidth $b_f$. Constraint (5) considers ARROW's restorable links and ensures the sum of the bandwidth of all restorable tunnels routed on link $e$ does not exceed $e$'s bandwidth capacity with a slack variables $\Delta_e^{z,q}$. For the slack variable $\Delta_e^{z,q}$ for link $e$ under scenario $q$ and LotteryTicket $z$, we also set a bound $M^{z,q}$ for the sum of $\Delta_e^{z,q}$ to ensure the total slack is within a reasonable region in Constraint (6).

**Phase I output.** To select the winning LotteryTicket without solving an ILP, Phase I outputs a floating-point slack variable $\Delta_e^{z,q}$ that allows the bandwidth allocation of restorable tunnels on IP link $e$ to go beyond the link's restorable capacity $r_e^{z,q}$. After solving ARROW's Phase I formulation, we run a post-processing step to find the winning LotteryTickets. This is done by comparing $r_e^{z,q} + \Delta_e^{z,q}$ with all LotteryTickets. Then, for each failure scenario $q$, we select the LotteryTicket $z$ with the minimum $\sum_{e \in E} max(0, \Delta_e^{z,q})$.[5] $Z^*$ is a set of size $|Q|$ that contains the winning LotteryTickets for each failure scenario $q \in Q$. ARROW maps the restoration plan $Z^*$ into wavelengths' reconfiguration rules and installs them on ROADM config files.

**Phase II input parameters.** As shown in Table 3, Phase II's input parameters include the winning LotteryTicket's restorable bandwidth capacity (denoted by $r_e^{*,q}$) and restorable tunnels (denoted by $Y_f^{*,q}$). These parameters are calculated from ARROW's Phase I post-processing step.

**Phase II optimization goal and constraints.** Phase II uses the same optimization goal (maximizing total network throughput) as Phase I. Moreover, the Constraints (7-9) of Phase II are the same as Constraints (1-3) in Phase I. However, Constraints (10-11) use

---

[4]In our evaluations, we experiment with $\alpha$ = 0.2, 0.1, and 0.05.

[5]This technique is similar to the ReLU function commonly used in machine learning.

the $r_e^{*,q}$ and $Y_f^{*,q}$ of the winning LotteryTickets (selected by Phase I) for each failure scenario $q$.

**Phase II output.** The output of the Phase II consists of two parts: (1) the *total* bandwidth $b_f$ that flow $f$ is permitted to utilize (across all of its tunnels in $T_f$); and (2) the allocation of $b_f$ over flow $f$'s tunnels $T_f$, denoted by $a_{f,t}$. Arrow periodically computes the optimal values of bandwidth allocations, based on the current demand matrix. Similar to prior TE schemes, after each TE run, Arrow finds the traffic splitting ratio for flow $f$ among its tunnels, calculated as $\omega_{f,t} = a_{f,t}/\sum_{t \in T_f} a_{f,t}$. These traffic splitting ratios are then installed on routers [53].[6]

**Probabilistic optimality guarantee.** Our goal is to find the optimal restoration plan $z^{opt}$ while maximizing throughput. However, as discussed in §3.1, simply formulating a joint IP/optical formulation is not scalable. As a result, Arrow relies on the LotteryTickets abstraction. However, given that the input LotteryTicket set $Z$ is generated by randomized rounding, Arrow's optimality depends on whether the optimal restoration plan appears in LotteryTickets set $Z$. Consequently, Arrow has a probabilistic optimality guarantee that depends on the probability of selecting $z^{opt}$ during its randomized rounding process. Note that we do not provide a deterministic optimality guarantee, hence, Arrow's solution can be sub-optimal. We leave finding a restoration-aware TE formulation, with a practical runtime and a deterministic optimality guarantee, to future work.

**Theorem** 3.1 (Arrow Probabilistic Optimality). *An* Arrow *TE with $|Z^q|$ LotteryTickets finds the optimal allocation under failure scenario $q$ with probability*

$$\rho^q = 1 - (1 - \kappa)^{|Z^q|} \tag{12}$$

*where $\kappa$ is the probability of finding the optimal LotteryTicket $z^{opt}$ for scenario $q$. This probability can be calculated as*

$$\kappa = \prod_{1 \le e \le n} \frac{1}{\delta} \times Pr\{round\ up/down\} \tag{13}$$

*where $1 \le e \le n$ is the index of the failed IP link under scenario $q$ (as defined in line 4 of Algorithm 1), and $Pr\{round\ up/down\}$ is the probability of rounding up or down to determine each failed IP link's restoration option (as defined in lines 6–11 of Algorithm 1).*

Please see Appendix A.3 for proof.

# 4 NOISE LOADING IN ARROW

For decades, wavelength reconfiguration in large-scale WANs has been deemed slow and complicated, even with ROADMs being already deployed. This is because newly reconfigured wavelengths change the power distribution over the fiber spectrum. Hence, the cascaded optical amplifiers along the fiber path need to adjust accordingly. This process introduces a non-trivial challenge to adjust amplifier gain configurations to equalize wavelength power/Signal-to-Noise Ratio (SNR) with repetitive observe-analyze-act loops which take a few minutes per amplifier to converge (Appendix A.7). This problem is commonly referred to as wavelength channel equalization and has been studied extensively in the optics community [7, 8, 16, 51, 55, 62, 67, 92]. Recent work used machine learning

---

[6]To avoid division by zero when the denominator is zero, in our code, we change the tunnels with $a_{f,t} = 0$ to $a_{f,t} = \epsilon$, where $\epsilon$ is a small number (i.e., $10^{-4}$).
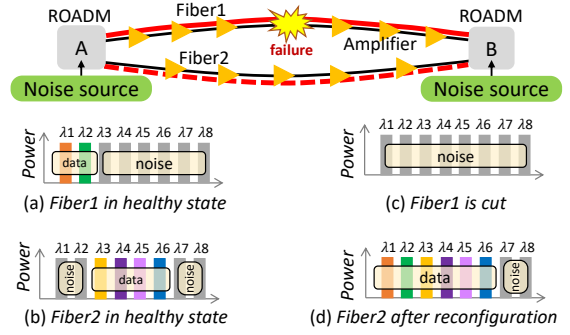


**Figure 9: Arrow's noise loading in action.**

to shorten this amplifier reconfiguration delay [90], but it requires a simulation-based firewall to ensure the parameters are safe to deploy on live production networks.

**Noise loading enables agile wavelength reconfigurations.** Arrow *bypasses* the amplifier reconfiguration latency by leveraging a technology called "noise loading" in modern optical backbones [33, 35, 83]. A programmable optical noise generation device, called the Amplified Spontaneous Emission (ASE) device [83], generates an optical noise signal for all *unused* wavelengths on all fibers. As a result, from the amplifier's perspective, all wavelengths are present at all times. Although all wavelength channels are turned up in all fibers, some wavelengths carry noise generated by the ASE device, and some carry IP layer data generated by routers.

**Noise loading example.** Consider the WAN depicted in Fig. 9. Assume all fibers can carry eight wavelengths in total.[7] In healthy state, two of Fiber1's available frequencies, $\lambda 1$ and $\lambda 2$, are connected to router ports and carry data while the other six frequencies, $\lambda 3$ to $\lambda 8$, are loaded with noise (Fig. 9(a)). Fiber2 has four wavelengths carrying data ($\lambda 3$, $\lambda 4$, $\lambda 5$, and $\lambda 6$), while the other frequencies ($\lambda 1$, $\lambda 2$, $\lambda 7$, and $\lambda 8$) are carrying noise (Fig. 9(b)). Now, assume Fiber1 is cut, causing wavelengths $\lambda 1$ and $\lambda 2$ to fail. Arrow reconfigures these two wavelengths onto Fiber2's $\lambda 1$ and $\lambda 2$ slots, which are initially loaded with noise, by configuring ROADM A and B, as well as their noise sources (Fig. 9(d)).[8] All of this is hidden from the amplifiers on both fibers because replacing noise with data is performed locally on the ROADMs (Appendix A.6). As a result, Arrow circumvents amplifier reconfiguration latency because amplifiers will no longer experience power changes as the total number of powered up spectrum frequencies (i.e., eight frequencies) is the same the entire time.

# 5 PRODUCTION-LEVEL TESTBED

**Setup.** To faithfully evaluate Arrow in a real-world setting, we select a subset of Facebook's global WAN and separate it from production for experimental purposes. Fig. 10 shows our testbed setup with four ROADMs, 34 amplifiers, and over 2,160 km unidirectional fiber. In our testbed, all hardware devices and software (control, monitoring, failure detection and management, etc.) are identical

---

[7]In practice, today's fibers can carry 48-96 wavelengths in the C-band range depending on channel frequency spacing [28, 38].
[8]Here, reconfiguring to $\lambda 1$ and $\lambda 2$ does not need frequency tuning. If $\lambda 1$ and $\lambda 2$ on fiber 2 are already carrying data, Arrow tunes the frequency of $\lambda 1$ and $\lambda 2$ to $\lambda 7$ and $\lambda 8$ to avoid colliding with wavelengths carrying data.
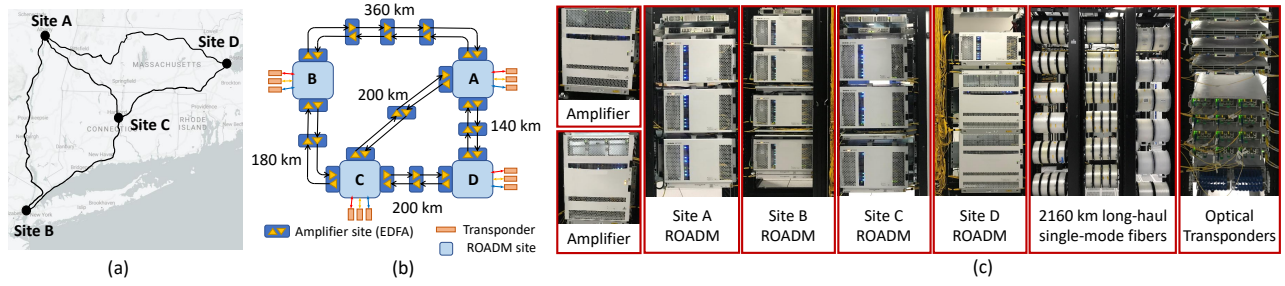
Figure 10: ARROW's production-level testbed. (a) Testbed topology emulating an optical backbone connecting 4 cites in North America. (b) Physical topology of the testbed. (c) Photo of the testbed.
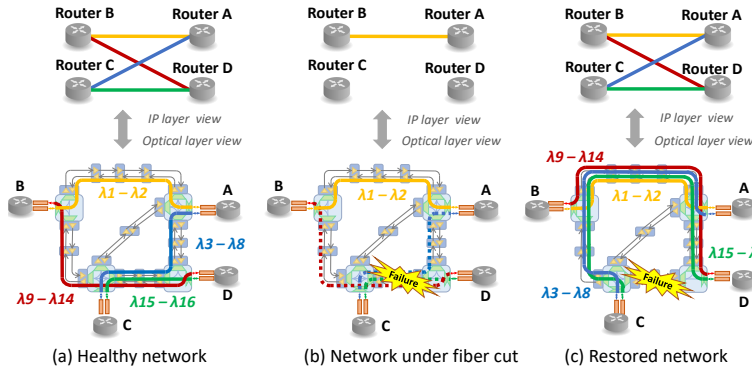


Figure 11: Restoring 2.8 Tbps of lost IP capacity with ARROW.



Figure 12: ARROW's restoration latency.

to Facebook's global WAN. This testbed demonstrates that ARROW can be readily deployed.

**A fiber cut restoration trial.** Fig. 11 shows an end-to-end experiment for a fiber cut restoration trial. Fig. 11(a) shows the network in healthy state. The top topology represents the IP-layer, and the bottom one represents optical fibers and devices (same layout as Fig. 10(b)). The color of each IP link matches its underlying fiber path. There are 16 wavelengths (each wavelength occupies 75 GHz spectrum frequency width modulated at 200 Gbps bandwidth capacity) grouped into 4 port-channels to support 4 IP links: A↔B (0.4 Tbps), A↔C (1.2 Tbps), B↔D (1.2 Tbps), and C↔D (0.4 Tbps). Fig. 11(b) shows the state of the network after fiber CD is cut. This fiber was carrying 14 wavelengths ($\lambda 3$-$\lambda 16$) and its cut caused 3 IP links, A↔C, B↔D, C↔D, to fail. To restore the IP-layer capacity, ARROW reconfigures the wavelengths as shown in Fig. 11(c).

**Quantifying the restoration latency.** Fig. 12 compares the restoration latency of ARROW with the state-of-the-art method [73]. We first use the current amplifier reconfiguration approach in Facebook to reconfigure the wavelengths (the same failure scenario as in Fig. 11). Figs. 12(a) and (b) show the normalized IP layer capacity and optical power measured on fiber AB; as the figure shows, it requires 1,021 seconds (≈17 mins) to restore 2.8 Tbps. We then use ARROW's noise loading device to bypass the amplifier reconfiguration latency. Figs. 12(c) and (d) show the entire restoration latency is eight seconds; i.e., 127× faster than the state-of-the-art method. We believe this latency can be further reduced to milliseconds with more advanced hardware, as shown in prior work [22]. Moreover, ARROW's reconfiguration does not affect existing wavelengths ($\lambda 1$ and $\lambda 2$) running on the fiber AB, as shown in Fig. 12(d).
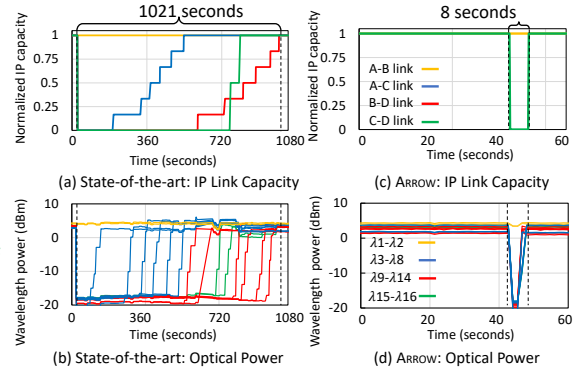
**Other factors affecting the latency.** Two other factors affect ARROW's restoration latency in practice: 1) wavelength tuning (if there is a frequency collision on the surrogate fiber path); 2) modulation change (if the length of the surrogate fiber path increases). Note that these two steps are optional, as appropriate, and can be adjusted in parallel with ROADM reconfiguration because they affect the transponders only. Prior work has demonstrated frequency tuning [10, 29] and modulation change [76, 77] in milliseconds.

## 6 LARGE-SCALE SIMULATIONS

We use simulations to quantify the performance gains of ARROW. Our simulation framework is implemented with the Julia programming language [15] and the Gurobi solver [66]. Our code is available online.[9] We compare ARROW to the following schemes:

• **FFC [63].** FFC is a failure-aware TE formulation that guarantees zero loss for up to $k$ IP-link failure scenarios. We extend FFC to the optical layer by considering scenarios with $k$ fiber cuts, and evaluate both $k = 1$ and $k = 2$ cases and refer to them as FFC-1 and FFC-2, respectively.

• **TeaVaR [17].** TeaVaR is also a failure-aware TE formulation but instead of absolute guarantees, it provides a probabilistic guarantee depending on the failure probability of fibers. In our simulations, we set TeaVaR's availability target ($\beta$) at 99.9%.

• **ECMP [21].** ECMP is not a failure-aware TE and serves as a baseline in our evaluations. It does not consider failure scenarios, hence it does not provide any guarantees with respect to failures; it simply assigns equal amount of traffic to all tunnels of each flow.
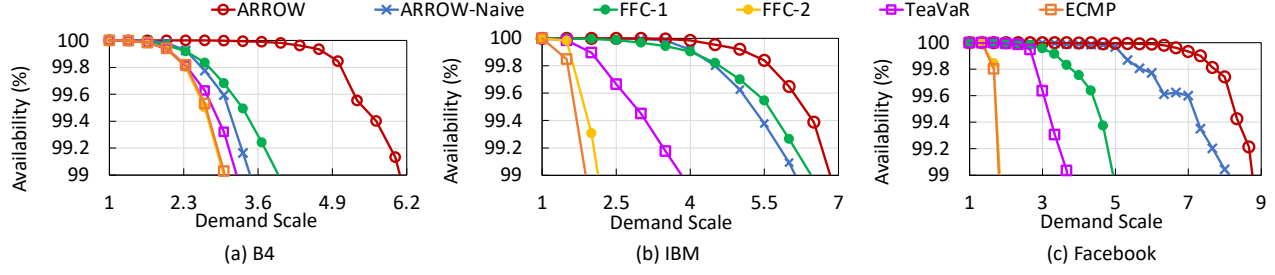
[9]http://arrow.csail.mit.edu

Figure 13: Availability vs. demand scales for ARROW and state-of-the-art failure-aware TE schemes.

Table 4: Network topologies used in our simulations.

| Topo. | # Routers / ROADMs | # Fibers | # IP links | # Traffic matrices |
|---|---|---|---|---|
| Facebook | 34/84 | 156 | 262 | 12 |
| IBM | 17/17 | 23 | 85 | 30 |
| B4 | 12/12 | 19 | 52 | 30 |

| | ARROW's gain in terms of satisfied demand | | | | |
|---|---|---|---|---|---|
| Availability | ARROW-Naive | FFC-1 | FFC-2 | TeaVaR | ECMP |
| 99.999% | 1.6× | 2.2× | 2.4× | 2.3× | 2.3× |
| 99.99% | 2.0× | 2.2× | 2.4× | 2.4× | 2.4× |
| 99.9% | 2.0× | 2.0× | 2.3× | 2.3× | 2.3× |
| 99% | 1.8× | 1.5× | 2.0× | 1.9× | 2.0× |

Table 5: ARROW's gain at different availability levels for B4 topology.

• **ARROW-Naive.** To evaluate the impact of ARROW's two-phased approach, we consider a naive version of ARROW, called ARROW-Naive, that consists of only phase II. Hence, instead of using LotteryTickets, this approach considers restoration solely at the optical layer without taking instantaneous traffic matrices into account. To do so, ARROW-Naive solves the RWA formulation (Appendix A.2) only once and uses its output as a winning LotteryTicket and bypasses Phase I at every TE run.

**Topologies.** We evaluate ARROW on three WAN topologies: Facebook, IBM, and B4 (Table 4). For Facebook topology, we use a subset of the optical-layer topology in production. For B4 and IBM, we take the topologies in [58] and use them as the optical-layer topology. Note that in large-scale WANs, the IP-layer topology tends to be denser than the optical-layer topology [65, 71]. To generate realistic IP-layer topologies, we measure the number of IP links per fiber and the number of wavelengths per IP link in Facebook (shown in Fig. 22 in Appendix A.8) and use these distributions to guide us to generate the IP-layer topologies. Unless otherwise stated, we use 120, 90, and 80 as the number of LotteryTickets for running ARROW on Facebook, IBM and B4, respectively.

**Traffic matrix.** For B4 and IBM networks, we use 30 traffic matrices from SMORE [58] generated by fitting the real-world traffic considering time variations and diurnal/weekly patterns. For the Facebook topology, we use 12 real traffic matrices from production.

**Tunnel selection.** ARROW is orthogonal to tunnel selection methods. In our evaluations, we use both fiber-disjoint routing and $k$-shortest path routing algorithms to route tunnels over the IP-layer topology, while ensuring that there is at least one residual tunnel for every flow under each failure scenario. We set the number of tunnels per flow at 8, 12, and 16, for B4, IBM, and Facebook, respectively.

**Fiber cut scenarios.** Following the methodology in TeaVaR [17], we use a Weibull distribution (shape=0.8, scale=0.02) to model the failure probability of each fiber. We then generate fiber cut scenarios using cutoff values of 0.001, 0.001, 0.0002 for B4, IBM, Facebook, respectively. Note that depending on the failure probabilities and the cutoff value, the generated fiber cut scenarios may contain both

single fiber cut or double fiber cut scenarios. When a fiber fails, all IP links on this fiber fail simultaneously.

**Demand scaling.** Given the fact that production WANs are over-provisioned, we start with a network state where 100% of traffic demand is satisfied. Similar to prior work [17, 48, 63], we scale the demand matrix uniformly to evaluate each TE's traffic allocation over tunnels under different traffic loads and failure scenarios.

## 6.1 Availability and Satisfied Demand Gains

In this section, we show that ARROW improves availability and throughput by restoring the lost IP capacity and reviving the IP-layer network.

**Availability metric definition.** Availability is a key metric to evaluate the satisfaction of Service Level Agreements (SLAs), and it is directly related to the revenue of network providers [41, 85]. Our availability metric is calculated as follows: for each topology and traffic matrix, we first solve each TE formulation to obtain traffic splitting rules. We then simulate all probabilistic failure scenarios and calculate the availability of each scenario based on the percentage of total demand satisfaction during that scenario. We then take the sum of the availabilities of all failure scenarios *weighted* by each scenario's probability as availability of a given traffic matrix. For each topology and demand scale, we then take the average availability across all traffic matrices.

**Impact of demand scaling on availability.** Fig. 13 shows the availability of different TE schemes on B4, IBM, and Facebook topologies. We focus on availability performance region larger than 99% because network operators need to maintain their network at high availability [41, 43, 85]. Fig. 13 shows that ARROW maintains higher availability levels as the demand is scaled for all three topologies. Specifically, we find that on B4 topology, ARROW can guarantee 99.99% availability, even when the traffic demand is scaled by 3.61×, while FFC-1 can sustain at most 1.63× demand increase at 99.99% availability. As a result, ARROW provides 2.2× gain in throughput compared to FFC-1 without sacrificing 99.99% availability. Table 5 summarizes ARROW's gains with respect to all considered prior
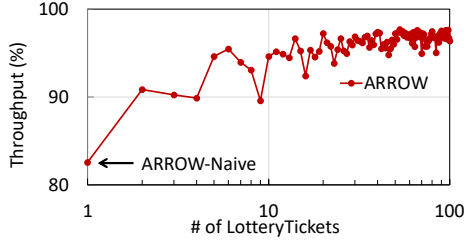
**Figure 14: Impact of number of LotteryTickets on Arrow's throughput in B4 topology.**



**Figure 15: Runtime of Arrow optimization.**

approaches for the B4 topology at different availability levels. We observe a similar trend for IBM and Facebook topologies, as shown in Fig. 13(b) and (c). At 99.99% availability, Arrow improves the network throughput by 1.6× and 2.4× compared to FFC-1. Note that FFC-1 only provides failure guarantees for single fiber cut scenarios while Arrow considers a combination of single and double fiber cuts. FFC-2 considers all double fiber failures but it has a considerably lower availability than Arrow (ECMP and FFC-2 curves are often overlapping). Although TeaVaR and Arrow both consider the same set of failure scenarios, Arrow outperforms TeaVaR by 2.4×, 2.8×, and 2.7× at 99.99% availability in B4, IBM, and Facebook, respectively.

## 6.2 Impact of LotteryTickets

**Throughput metric definition.** Network throughput is another core metric to evaluate TE algorithms because it shows the total traffic that a network can accommodate. For each topology and traffic matrix, our throughput metric is calculated as the ratio of total admissible bandwidth over total demand ($\frac{\sum_f b_f}{\sum_f d_f}$) returned by the TE optimization formulation. We then take the average network throughput across all traffic matrices.

**Impact of number of LotteryTickets on throughput.** Fig. 14 shows the impact of the number of LotteryTickets on Arrow's network throughput for B4 topology when the demand is scaled by 4.2×. The figure shows that when the number of LotteryTickets is small, the throughput fluctuates. This is because LotteryTickets are generated using randomized rounding, hence, more LotteryTickets are probabilistically better. When the number of LotteryTickets is one, it means we only have one restoration candidate for each failure scenario and hence it represents the Arrow-Naive approach where the restoration plan comes from solving the optical restoration RWA formulation (Appendix A.2) offline. As the number of LotteryTickets increases, Arrow's throughput gradually increases with less fluctuations until it reaches a plateau reflecting that the LotteryTickets have already covered a good set of restoration candidates, and continuing to add new LotteryTickets does not help much. To find a balance between TE execution time and throughput, the operator should select the appropriate number of LotteryTickets.

**TE optimization runtime.** We now compare Arrow's optimization runtime for different number of LotteryTickets. Arrow's optimization is formulated as an LP, and is solvable in polynomial time using Gurobi [66]. Fig. 15 presents Gurobi's solve time for Arrow TE optimization (Phase I + Phase II runtime) on a Linux server with AMD EPYC 7502P 32-Core CPU processor and 256 GB
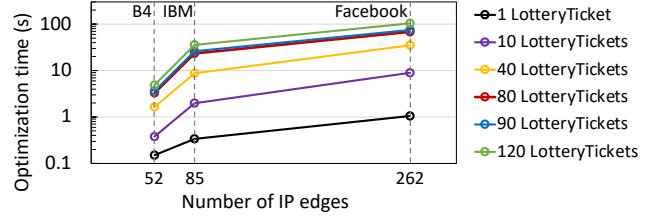
RAM. Note that this runtime only captures the *optimization solve time* and excludes the time it takes to build the optimization model. The figure shows that, Arrow's runtime increases as the number of LotteryTickets increase. For the Facebook topology with 120 LotteryTickets, Arrow's formulation is solved within 104 seconds. As mentioned in Section 3.1, the deadline for our TE runtime is 5 minute. Hence, Arrow is within the acceptable runtime range in Facebook.

## 6.3 Cost Savings

**Availability-guaranteed throughput.** To compare the cost associated with each TE scheme, following TeaVaR's approach, we first compute the availability-guaranteed throughput for each TE. This is because different TE algorithms that we consider provide different availability guarantees. For instance, FFC-1 guarantees 100% availability for all single fiber cut scenarios, but it does not guarantee anything for double fiber cuts. On the other hand, FFC-2 guarantees 100% availability for all double fiber cuts and hence achieves lower throughput. To make apples-with-apples comparison, we calculate the availability-guaranteed throughput. Specifically, for a given availability target (e.g., $\beta = 99.9\%$), we iterate over all failure scenarios of interest to compute the normalized demand loss for each scenario. We then sort the failure scenarios based on their loss values and find the scenario at the $\beta$-percentile. The normalized satisfied demand (1 - loss) of this scenario is reported as the availability-guaranteed throughput. In other words, the throughput is guaranteed to be no smaller than this value for $\beta$ percent of failure scenarios.

**Number of required router ports.** The number of required router ports to sustain a highly available network directly relates to the cost of the network. To calculate the number of required router ports, we find the worst-case traffic allocated on each IP link $e$ across all failure scenarios ($CAP_e$). We then calculate $CAP = \sum_e CAP_e$ to find the required network capacity for the entire topology. To make a fair comparison across different TE schemes, we then normalize $CAP$ by the availability-guaranteed throughput value as a proxy for the number of required router ports for each TE scheme and network topology. Fig. 16 shows the number of required router ports to achieve the same availability-guaranteed throughput with $\beta = 99.9\%$ *availability target* for different TE algorithms. To put Arrow's savings into perspective, we also calculate the minimum number of required router ports by considering a *hypothetical TE* that can achieve 100% availability at all times by fully restoring every failure scenario. This reflects a TE that does not require any *over-provisioning of router ports* to achieve 100% availability. We call this approach *Fully Restorable TE* and use it as the baseline in Fig. 16. The figure shows that Arrow has a fundamental advantage over
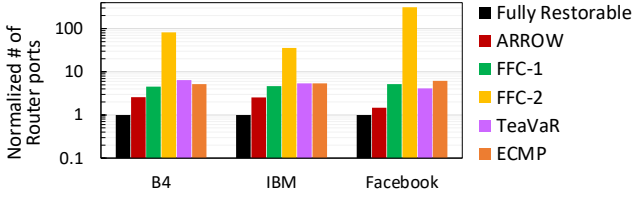
**Figure 16: Number of router ports required for different TE algorithms to support the same throughput with the same availability level (99.9%).**

existing TEs because it can restore lost IP capacity. This feature allows ARROW to provide the same level of availability with less over-provisioning. Specifically, for the Facebook topology, ARROW requires 2.8× fewer router ports than the best failure-aware TE (i.e., TeaVaR). Importantly, TeaVaR, FFC-1, and FFC-2 require 4.1×, 5.2×, 311.4× more router ports compared to the fully restorable TE. However, ARROW only requires 1.5× more router ports to sustain high availability even with partially restorable fibers (see §2.3). Although our fully restorable TE case is a hypothetical best case scenario, this result highlights ARROW's ability to maintain high throughput and high availability under failures without requiring extensive over-provisioning.

**Number of required transponders.** Similar to router ports, the number of required transponders also impact cost. In general, there is a 1-to-1 mapping between router ports and transponders (Fig. 1). Hence, the cost savings for router ports directly relate to cost savings for transponders as well.

## 7 RELATED WORK

**Traffic Engineering.** Traffic engineering is an important topic in WANs [6, 17, 42, 43, 47, 48, 57, 58, 63, 79]. Work related to ARROW includes failure-aware TE techniques [17, 48, 63], where the TE formulation considers failure scenarios and pre-allocates enough headroom on links so that when failures happen, traffic loss is minimized (or is zero). Although such techniques embed failure-recovery constraints in the TE formulation, in the case of fiber cuts, they end up under-utilizing the WAN significantly. TeaVaR [17] improves the utilization by assigning a probability to each failure scenario but it still needs to allocate headroom for probable failures. The first row in Table 10 (see Appendix A.9) illustrates the properties of this class of solutions. There are also other failure-oblivious TE algorithms that aim at assigning traffic that respecting link capacity [42], distributing traffic to equalize link utilization with traffic-oblivious tunnels [58], or contracting network topologies for runtime optimization [6]. These algorithms generally do not consider failure scenarios in their formulation and can only respond to failures in a reactive way without performance guarantees.

**Optical path protection.** Optical path protection techniques [14, 19, 59, 68, 80, 81, 84] pre-allocate failover paths solely on the optical domain using a device called the Optical Transport Network (OTN) switch. This is done by statically assigning a set of standby failover paths to each fiber during the capacity planning phase. When a failure happens, the OTN device quickly shifts the traffic from the failed fiber to its active back up path without notifying the TE. The second row in Table 10 illustrates this class of solutions. Although

this approach saves on router ports, it still needs to pre-allocate transponders and keeps them idle to be prepared for fast failover. Moreover, since the failover is entirely configured in the optical layer, the TE is blind to the extra available capacity and cannot utilize it optimally.

**Classical optical restoration.** Classical optical restoration techniques are the most relevant work to this paper. Although the benefits of optical restoration have been demonstrated in prior work [30, 32, 45, 49, 52, 56, 74, 82], to the best of our knowledge, there is no practical study of optical restoration in modern WANs. ARROW makes two novel contributions. First, ARROW augments today's TE formulations to capture partial restoration candidates for IP links (§3). Second, ARROW uses a noise source to fully populate the amplifiers' spectrum to bypass their reconfiguration time (§4), achieving an end-to-end failover latency of eight seconds on a WAN-scale testbed (§5).

**Reconfigurable WANs.** Recently, there have been several proposals to enable reconfigurable WANs [18, 27, 33, 34, 37, 49, 71, 77, 87, 89, 91]. Iris proposed an all-optical circuit switched network to interconnect datacenter sites that are only a few tens of kilometers apart (metro-level) [33]. In contrast, ARROW considers sites that are thousands of kilometers apart and focuses on fiber cut restoration. Another class of prior work proposed enabling reconfigurability in the optical domain to accommodate traffic matrix changes [18, 27, 34, 49, 71, 87, 91]. Notably, OWAN demonstrated reconfiguring optical wavelengths to adapt topology to achieve better bulk transfer performance in the WAN [49]. However, OWAN did not consider failures and its emulated ROADMs did not consider reconfiguration latency. RADWAN [77] proposed changing transponder modulations according to changes to OSNR on fiber paths to achieve better link utilization and availability. But RADWAN did not consider fiber cuts. When a fiber cut happens, RADWAN cannot change the modulation of the wavelengths because the fiber is down and changing the modulations will not help. OptFlow [37] proposed a graph model to enable optical reconfigurability with no change to the TE. But it did not consider partial restoration; hence, the approach is not portable to ARROW.

## 8 CONCLUSION

We propose a restoration-aware TE system, called ARROW, to proactively consider partially restorable failures when optimizing traffic allocations. While the restoration is done on the optical layer and TE is done on the IP layer, we avoid the computational complexity of conventional cross-layer formulations by designing a novel abstraction, called LotteryTicket, to feed only essential information into the TE formulation to meet the stringent TE runtime requirement. Our experiments show ARROW supports up to 2.0×–2.4× more demand without compromising availability at 99.99% availability. This work does not raise any ethical issues.

# REFERENCES

[1] 2014. Infinera introduces flexible grid 500G super-channel ROADM. (March 2014). http://www.gazettabyte.com/home/2014/3/14/infinera-introduces-flexible-grid-500g-super-channel-roadm.html.

[2] 2020. Automatic and Manual optical redundant failover switch. (2020). http://www.comlaninc.com/products/fiber-optic-products/id/23/cl-fos.

[3] 2021. Adva ROADM. (2021). https://www.adva.com/en/products/technology/roadm.

[4] 2021. CDC ROADM Applications and Cost Comparison. (2021). https://www.ofcconference.org/getattachment/188d14da-88ba-4a63-91d6-1cc14b335d8b/CDC-ROADM-Applications-and-Cost-Comparison.aspx.

[5] 2021. Optical Transceivers. (2021). https://www.smartoptics.com/article/optical-transceivers-turning-data-into-light/.

[6] Firas Abuzaid, Srikanth Kandula, Behnaz Arzani, Ishai Menache, Matei Zaharia, and Peter Bailis. 2021. Contracting Wide-area Network Topologies to Solve Flow Problems Quickly. In *18th USENIX Symposium on Networked Systems Design and Implementation (NSDI 21)*. USENIX Association, 175–200. https://www.usenix.org/conference/nsdi21/presentation/abuzaid.

[7] Choudhury A Al Sayeed, David C Bownass, David W Boertjes, and GAO Shiyu. 2017. Spectrum controller systems and methods in optical networks. (Feb. 21 2017). US Patent 9,577,763.

[8] Choudhury A Al Sayeed, Dave C Bownass, and Edward Chen. 2018. Systems and methods modeling optical sources in optical spectrum controllers for control thereof. (May 29 2018). US Patent 9,986,317.

[9] Mohammadreza Alizadeh Attar, Sha Ma, and Thomas J Edsall. 2017. Randomized per-packet port channel load balancing. (March 7 2017). US Patent 9,590,914.

[10] Hitesh Ballani, Paolo Costa, Raphael Behrendt, Daniel Cletheroe, István Haller, Krzysztof Jozwik, Fotini Karinou, Sophie Lange, Kai Shi, Benn Thomsen, and Hugh Williams. 2020. Sirius: A Flat Datacenter Network with Nanosecond Optical Switching. In *SIGCOMM'20*. 782–797.

[11] D. Banerjee and B. Mukherjee. 1996. A practical approach for routing and wavelength assignment in large wavelength-routed optical networks. *IEEE Journal on Selected Areas in Communications* 14, 5 (1996), 903–908. https://doi.org/10.1109/49.510913

[12] Dhritiman Banerjee and Biswanath Mukherjee. 2000. Wavelength-routed optical networks: Linear formulation, resource budgeting tradeoffs, and a reconfiguration study. *IEEE/ACM Transactions on networking* 8, 5 (2000), 598–607.

[13] Ajay Kumar Bangla, Alireza Ghaffarkhah, Ben Preskill, Bikash Koley, Christopher Albrecht, Emilie Danna, Joe Jiang, and Xiaoxue Zhao. 2015. Capacity planning for the Google backbone network. In *ISMP (2015)*.

[14] Marco Bertolini, Olivier Rocher, Arnaud Bisson, Pascal Pecci, and Giovanni Bellotti. 2012. Benefits of OTN switching introduction in 100Gb/s optical transport networks. In *OFC/NFOEC*. IEEE, 1–3.

[15] Jeff Bezanson, Alan Edelman, Stefan Karpinski, and Viral B Shah. 2017. Julia: A fresh approach to numerical computing. *SIAM review* 59, 1 (2017), 65–98.

[16] Berk Birand, Howard Wang, Keren Bergman, Dan Kilper, Thyaga Nandagopal, and Gil Zussman. 2014. Real-time power control for dynamic optical networks: Algorithms and experimentation. *IEEE Journal on Selected Areas in Communications* 32, 8 (2014), 1615–1628.

[17] Jeremy Bogle, Nikhil Bhatia, Manya Ghobadi, Ishai Menache, Nikolaj Bjørner, Asaf Valadarsky, and Michael Schapira. 2019. TEAVAR: Striking the Right Utilization-Availability Balance in WAN Traffic Engineering. In *Proceedings of the ACM Special Interest Group on Data Communication (SIGCOMM'19)*. Association for Computing Machinery, New York, NY, USA, 29–43. https://doi.org/10.1145/3341302.3342069

[18] Andrew Brzezinski and Eytan Modiano. 2005. Dynamic reconfiguration and routing algorithms for IP-over-WDM networks with stochastic traffic. In *INFOCOM*.

[19] D. Cavendish. 2000. Evolution of optical transport technologies: from SONET/SDH to WDM. *IEEE Communications Magazine* 38, 6 (2000), 164–172. https://doi.org/10.1109/35.846090

[20] Bowen Chen, Jie Zhang, Yongli Zhao, Jason P. Jue, Jinyan Liu, Shanguo Huang, and Wanyi Gu. 2014. Spectrum block consumption for shared-path protection with joint failure probability in flexible bandwidth optical networks. *Optical Switching and Networking* 13 (2014), 49 – 62. https://doi.org/10.1016/j.osn.2014.01.001

[21] Marco Chiesa, Guy Kindler, and Michael Schapira. 2016. Traffic engineering with equal-cost-multipath: An algorithmic perspective. *IEEE/ACM Transactions on Networking* 25, 2 (2016), 779–792.

[22] Angela L Chiu, Gagan Choudhury, George Clapp, Robert Doverspike, Mark Feuer, Joel W Gannett, Janet Jackel, Gi Tae Kim, John G Klincewicz, Taek Jin Kwon, et al. 2011. Architectures and protocols for capacity efficient, highly dynamic and highly resilient core networks. *IEEE/OSA Journal of Optical Communications and Networking* 4, 1 (2011), 1–14.

[23] Angela L Chiu, Gagan Choudhury, Mark D Feuer, John L Strand, and Sheryl L Woodward. 2011. Integrated restoration for next-generation IP-over-optical networks. *Journal of Lightwave Technology* 29, 6 (2011), 916–924.

[24] Junho Cho and Peter J Winzer. 2019. Probabilistic constellation shaping for optical fiber communications. *Journal of Lightwave Technology* 37, 6 (2019), 1590–1607.

[25] Cisco. 2019. Chapter: Configuring Port Channels. (2019). Cisco Nexus 5000 Series NX-OS Software Configuration Guide.

[26] Cisco. 2020. Cisco Network Convergence System 1004 L-Band Transponder Line Card Data Sheet . https://www.cisco.com/c/en/us/products/collateral/optical-networking/network-convergence-system-1000-series/datasheet-c78-743956.html. (2020).

[27] J. Cox. 2015. SDN control of a coherent Open Line System. In *2015 Optical Fiber Communications Conference and Exhibition (OFC)*. 1–1. https://doi.org/10.1364/OFC.2015.M3H.4

[28] O Gonzalez de Dios, R Casellas, F Zhang, X Fu, D Ceccarelli, and I Hussain. 2015. Framework and requirements for GMPLS-based control of flexi-grid dense wavelength division multiplexing (DWDM) networks. In *IETF RFC 7698*.

[29] S. Dhoore, G. Roelkens, and G. Morthier. 2019. Fast Wavelength-Tunable Lasers on Silicon. *IEEE Journal of Selected Topics in Quantum Electronics* 25, 6 (2019), 1–8. https://doi.org/10.1109/JSTQE.2019.2912034

[30] Bharat T Doshi, Subrahmanyam Dravida, P Harshavardhana, Oded Hauser, and Yufei Wang. 1999. Optical network design and restoration. *Bell Labs Technical Journal* 4, 1 (1999), 58–84.

[31] Robert Doverspike and Jennifer Yates. 2001. Challenges for MPLS in optical network restoration. *IEEE Communications magazine* 39, 2 (2001), 89–96.

[32] Robert D Doverspike. 2020. Carrier Network Architectures and Resiliency. In *Springer Handbook of Optical Networks*. Springer, 399–446.

[33] Vojislav Dukic, Ginni Khanna, Christos Gkantsidis, Thomas Karagiannis, Francesca Parmigiani, Ankit Singla, Mark Filer, Jeffrey L. Cox, Anna Ptasznik, Nick Harland, Winston Saunders, and Christian Belady. 2020. Beyond the Mega-Data Center: Networking Multi-Data Center Regions. In *Proceedings of the Annual Conference of the ACM Special Interest Group on Data Communication on the Applications, Technologies, Architectures, and Protocols for Computer Communication (SIGCOMM '20)*. Association for Computing Machinery, New York, NY, USA, 765–781. https://doi.org/10.1145/3387514.3406220

[34] Mark Filer, Jamie Gaudette, Monia Ghobadi, Ratul Mahajan, Tom Issenhuth, Buddy Klinkers, and Jeff Cox. 2016. Elastic Optical Networking in the Microsoft Cloud. *Journal of Optical Communications and Networking* 8, 7, A45–A54.

[35] Mark Filer, Jamie Gaudette, Yawei Yin, Denizcan Billor, Zahra Bakhtiari, and Jeffrey L Cox. 2019. Low-margin optical networking at cloud scale. *IEEE/OSA Journal of Optical Communications and Networking* 11, 10 (2019), C94–C108.

[36] Finisar. 2020. Finisar Dual Wavelength Selective Switch (WSS) . https://finisarwss.com/wp-content/uploads/2020/07/FinisarWSS_Dual_Wavelength_Selective_Switch_ProductBrief_Jul2020.pdf. (2020).

[37] Klaus-Tycho Foerster, Long Luo, and Manya Ghobadi. 2020. OptFlow: A Flow-Based Abstraction for Programmable Topologies. In *Proceedings of the Symposium on SDN Research (SOSR '20)*. Association for Computing Machinery, New York, NY, USA, 96–102. https://doi.org/10.1145/3373360.3380840

[38] ITU-T Recommendation G.694.1. 2012. Spectral grids for WDM applications: DWDM frequency grid. (2012). https://www.itu.int/rec/T-REC-G.694.1/en.

[39] Ori Gerstel, Clarence Filsfils, Thomas Telkamp, Matthias Gunkel, Martin Horneffer, Victor Lopez, and Arturo Mayoral. 2014. Multi-layer capacity planning for IP-optical networks. *IEEE Communications Magazine* 52, 1 (2014), 44–51.

[40] Jennifer Gossels, Gagan Choudhury, and Jennifer Rexford. 2019. Robust network design for IP/optical backbones. *J. Opt. Commun. Netw.* 11, 8 (Aug 2019), 478–490. https://doi.org/10.1364/JOCN.11.000478

[41] Tamás Hauer, Philipp Hoffmann, John Lunney, Dan Ardelean, and Amer Diwan. 2020. Meaningful availability. In *17th {USENIX} Symposium on Networked Systems Design and Implementation ({NSDI} 20)*. 545–557.

[42] 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'13* (2013), 12.

[43] Chi-Yao Hong, Subhasree Mandal, Mohammad Al-Fares, Min Zhu, Richard Alimi, Chandan Bhagat, Sourabh Jain, Jay Kaimal, Shiyu Liang, Kirill Mendelev, et al. 2018. B4 and after: managing hierarchy, partitioning, and asymmetry for availability and scale in google's software-defined WAN. In *SIGCOMM'18*. 74–87.

[44] IEEE. [n. d.]. IEEE 802.3ad Link Aggregation. https://www.ieee802.org/3/ad/. ([n. d.]).

[45] Rainer R Iraschko and Wayne D Grover. 2000. A highly efficient path-restoration protocol for management of optical network transport integrity. *IEEE Journal on Selected Areas in Communications* 18, 5 (2000), 779–794.

[46] Khalil A Jabr, Sudhakar Shenoy, and Dileep K Devireddy. 2011. Distribution of Packets Among PortChannel Groups of PortChannel Links. (May 12 2011). US Patent App. 12/645,564.

[47] 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* (2013), 12.

[48] Chuan Jiang, Sanjay Rao, and Mohit Tawarmalani. 2020. PCF: Provably Resilient Flexible Routing. In *Proceedings of the Annual Conference of the ACM*

*Special Interest Group on Data Communication on the Applications, Technologies, Architectures, and Protocols for Computer Communication (SIGCOMM '20).* Association for Computing Machinery, New York, NY, USA, 139–153. https://doi.org/10.1145/3387514.3405858

[49] 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 *Proceedings of the 2016 ACM SIGCOMM Conference.* 87–100.

[50] Masahiko Jinno, Bartlomiej Kozicki, Hidehiko Takara, Atsushi Watanabe, Yoshiaki Sone, Takafumi Tanaka, and Akira Hirano. 2010. Distance-adaptive spectrum resource allocation in spectrum-sliced elastic optical path network [topics in optical communications]. *IEEE Communications Magazine* 48, 8 (2010), 138–145.

[51] Joseph Junio, Daniel C Kilper, and Vincent WS Chan. 2012. Channel power excursions from single-step channel provisioning. *IEEE/OSA Journal of Optical Communications and Networking* 4, 9 (2012), A1–A7.

[52] Takuya Kanai, Yumiko Senoo, Kota Asaka, Jun Sugawa, Hideaki Tamai, Hiroyuki Saito, Naoki Minato, Atsushi Oguri, Seiya Sumita, Takehiro Sato, et al. 2018. Novel automatic service restoration technique by using self-reconfiguration of network resources for a disaster-struck metro-access network. *Journal of Lightwave Technology* 36, 8 (2018), 1516–1523.

[53] Nanxi Kang, Monia Ghobadi, John Reumann, Alexander Shraer, and Jennifer Rexford. 2015. Efficient traffic splitting on commodity switches. In *Proceedings of the 11th ACM Conference on Emerging Networking Experiments and Technologies.* 1–13.

[54] Ginni Khanna, Shengxiang Zhu, Mark Filer, Christos Gkantsidis, Francesca Parmigiani, and Thomas Karagiannis. 2020. Towards all optical DCI networks, In Optical Fiber Communication Conference (OFC) 2020. *Optical Fiber Communication Conference (OFC) 2020*, W2A.33. https://doi.org/10.1364/OFC.2020.W2A.33

[55] D Kilper, M Bhopalwala, H Rastegarfar, and W Mo. 2015. Optical power dynamics in wavelength layer software defined networking. In *Photonic Networks and Devices.* Optical Society of America, NeT2F–2.

[56] Adil Kodian and Wayne D Grover. 2005. Failure-independent path-protecting p-cycles: Efficient and simple fully preconnected optical-path protection. *Journal of lightwave technology* 23, 10 (2005), 3241.

[57] Praveen Kumar, Chris Yu, Yang Yuan, Nate Foster, Robert Kleinberg, and Robert Soulé. 2018. YATES: Rapid prototyping for traffic engineering systems. In *Proceedings of the Symposium on SDN Research.* 1–7.

[58] Praveen Kumar, Yang Yuan, Chris Yu, Nate Foster, Robert Kleinberg, Petr Lapukhov, Chiun Lin Lim, and Robert Soulé. 2018. Semi-oblivious traffic engineering: The road not taken. In *NSDI'18.* 157–170.

[59] Danny Lahav, Haim Moshe, Arie Menscher, and Aryeh Lezerovitz. 2006. Combined SONET/SDH and OTN architecture. (Sept. 12 2006). US Patent 7,106,968.

[60] Youngseok Lee, Yongho Seok, Yanghee Choi, and Changhoon Kim. 2002. A constrained multipath traffic engineering scheme for MPLS networks. In *ICC (2002).* IEEE, 2431–2436.

[61] George Leopold. 2017. Building Express Backbone: Facebook's new long-haul network. http://code.facebook.com/posts/1782709872057497/. (2017).

[62] Yao Li and Daniel C Kilper. 2018. Optical physical layer SDN. *Journal of Optical Communications and Networking* 10, 1 (2018), A110–A121.

[63] Hongqiang Harry Liu, Srikanth Kandula, Ratul Mahajan, Ming Zhang, and David Gelernter. 2014. Traffic engineering with forward fault correction. In *Proceedings of the 2014 ACM conference on SIGCOMM.* 527–538.

[64] Biswanath Mukherjee. 2006. *Optical WDM networks.* Springer Science & Business Media.

[65] Biswanath Mukherjee, Dhritiman Banerjee, Sav Ramamurthy, and Amarnath Mukherjee. 1996. Some principles for designing a wide-area WDM optical network. *IEEE/ACM transactions on networking* 4, 5 (1996), 684–696.

[66] Gurobi Optimization. 2021. Gurobi optimizer. (2021). http://www.gurobi.com.

[67] Yan Pan, Daniel C Kilper, Annalisa Morea, Joseph Junio, and Vincent WS Chan. 2012. Channel power excursions in GMPLS end-to-end optical restoration with single-step wavelength tuning. In *OFC/NFOEC.* IEEE, 1–3.

[68] Coriant White Paper. [n. d.]. The role of OTN switching in 100G & beyond transport networks. ([n. d.]). https://www.ofcconference.org/getattachment/90c0e6a4-08c1-45fb-a7f2-2957d444dc7d/The-Role-of-OTN-Switching-in-100G-Beyond-Transpo.aspx.

[69] Abhinav Pathak, Ming Zhang, Y Charlie Hu, Ratul Mahajan, and Dave Maltz. 2011. Latency inflation with MPLS-based traffic engineering. In *Proceedings of the 2011 ACM SIGCOMM conference on Internet measurement conference.* 463–472.

[70] Prabhakar Raghavan and Clark D Tompson. 1987. Randomized rounding: a technique for provably good algorithms and algorithmic proofs. *Combinatorica* 7, 4 (1987), 365–374.

[71] Byrav Ramamurthy and Ashok Ramakrishnan. 2000. Virtual topology reconfiguration of wavelength-routed optical WDM networks. In *GLOBECOM.*

[72] S Ramamurthy and Biswanath Mukherjee. 1999. Survivable WDM mesh networks. II. Restoration. In *1999 IEEE International Conference on Communications (Cat. No. 99CH36311)*, Vol. 3. IEEE, 2023–2030.

[73] Kim B Roberts, James Harley, and David Boertjes. 2020. Adjustment of control parameters of section of optical fiber network. (Sept. 22 2020). US Patent 10,784,980.

[74] Ali Najib Saleh, Haig Michael Zadikian, Zareh Baghdasarian, and Vahid Parsi. 2005. Method of reducing traffic during path restoration. (Feb. 1 2005). US Patent 6,850,486.

[75] Nicola Sambo, Alessio Ferrari, Antonio Napoli, Nelson Costa, João Pedro, Bernd Sommerkorn-Krombholz, Piero Castoldi, and Vittorio Curri. 2020. Provisioning in Multi-Band Optical Networks. *Journal of Lightwave Technology* 38, 9 (2020), 2598–2605.

[76] Rachee Singh, Monia Ghobadi, Klaus-Tycho Foerster, Mark Filer, and Phillipa Gill. 2017. Run, walk, crawl: Towards dynamic link capacities. In *Proceedings of the 16th ACM Workshop on Hot Topics in Networks.* 143–149.

[77] Rachee Singh, Manya Ghobadi, Klaus-Tycho Foerster, Mark Filer, and Phillipa Gill. 2018. RADWAN: Rate Adaptive Wide Area Network. In *Proceedings of the 2018 Conference of the ACM Special Interest Group on Data Communication (SIGCOMM '18).* ACM, New York, NY, USA, 547–560. https://doi.org/10.1145/3230543.3230570

[78] Andrew D. Smith. 2003. Probabilistic Methods in Integer Programming. http://smithlabresearch.org/downloads/randomized_rounding_and_integer_programming.pdf. (2003).

[79] Martin Suchara, Dahai Xu, Robert Doverspike, David Johnson, and Jennifer Rexford. 2011. Network Architecture for Joint Failure Recovery and Traffic Engineering. In *ACM SIGMETRICS (2011).*

[80] Massimo Tornatore, Guido Maier, and Achille Pattavina. 2005. Availability design of optical transport networks. *IEEE Journal on Selected Areas in Communications* 23, 8 (2005), 1520–1532.

[81] Jean-Philippe Vasseur, Mario Pickavet, and Piet Demeester. 2004. *Network Recovery: Protection and Restoration of Optical, SONET-SDH, IP, and MPLS.* Morgan Kaufmann Publishers Inc., San Francisco, CA, USA.

[82] Ann Von Lehmen, Robert Doverspike, George Clapp, Douglas M Freimuth, Joel Gannett, Aleksandar Kolarov, Haim Kobrinski, Christian Makaya, Emmanuil Mavrogiorgis, Jorge Pastor, et al. 2015. CORONET: Testbeds, demonstration, and lessons learned. *IEEE/OSA Journal of Optical Communications and Networking* 7, 3 (2015), A447–A458.

[83] Helen Xenos. 2020. How 3 technology advancements provide new options for scaling your optical network. (2020). https://www.ciena.com/insights/articles/how-3-technology-advancements-provide-new-options-for-scaling-your-optical-network.html.

[84] Tiejun J Xia, Steven Gringeri, and Masahito Tomizawa. 2012. High-capacity optical transport networks. *IEEE Communications Magazine* 50, 11 (2012), 170–178.

[85] Yiting Xia, Ying Zhang, Zhizhen Zhong, Guanqing Yan, Chiunlin Lim, Satyajeet Singh Ahuja, Soshant Bali, Alexander Nikolaidis, Kimia Ghobadi, and Manya Ghobadi. 2021. A Social Network Under Social Distancing: Risk-Driven Backbone Management During COVID-19 and Beyond.. In *NSDI.* 217–231.

[86] Jin Y Yen. 1971. Finding the k shortest loopless paths in a network. *management Science* 17, 11 (1971), 712–716.

[87] Y. Yoshida, A. Maruta, K.-I. Kitayama, M. Nishihara, T. Tanaka, T. Takahara, J.C. Rasmussen, N. Yoshikane, T. Tsuritani, I. Morita, Shuangyi Yan, Yi Shu, Yan Yan, R. Nejabati, G. Zervas, D. Simeonidou, R. Vilalta, R. Munoz, 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. *Lightwave Technology, Journal of* 33, 3 (Feb 2015), 609–617.

[88] Hui Zang, Jason P Jue, Biswanath Mukherjee, et al. 2000. A review of routing and wavelength assignment approaches for wavelength-routed optical WDM networks. *Optical networks magazine* 1, 1 (2000), 47–60.

[89] Chengliang Zhang, Junjie Li, Haiqiang Wang, Aihua Guo, and Christopher Janz. 2019. Evaluation of Dynamic Optical Service Restoration on a Large-Scale ROADM Mesh Network. *IEEE Communications Magazine* 57, 4 (2019), 138–143.

[90] Zhizhen Zhong, Manya Ghobadi, Maximilian Balandat, Sanjeevkumar Katti, Abbas Kazerouni, Jonathan Leach, Mark McKillop, and Ying Zhang. 2021. BOW: First Real-World Demonstration of a Bayesian Optimization System for Wavelength Reconfiguration. In *2021 Optical Fiber Communications Conference (OFC).*

[91] Zhizhen Zhong, Nan Hua, Massimo Tornatore, Jialong Li, Yanhe Li, Xiaoping Zheng, and Biswanath Mukherjee. 2019. Provisioning Short-Term Traffic Fluctuations in Elastic Optical Networks. *IEEE/ACM Transactions on Networking* 27, 4 (2019), 1460–1473.

[92] John Zyskind and Atul Srivastava. 2011. *Optically amplified WDM networks.* Academic press.

# A APPENDIX

Appendices are supporting material that has not been peer-reviewed.

## A.1 Path Inflation and Modulation Change After Restoration

**Fiber path length inflation.** An important factor in ARROW is the potential path inflation when the length of the surrogate restoration fiber path (R-path) is significantly longer than the original primary fiber path (P-path). Fig. 17 plots the CDF of the path inflation ratio (restoration path length divided by primary path length) for ARROW to perform optical restoration with and without frequency tuning at Facebook. Interestingly, the figure shows that, on average, 50% of IP links' restoration paths are *shorter* than the corresponding primary paths. This means the modulation formats of the restored wavelengths do not need to be reconfigured, thus simplifying AR-ROW's operations. For the remaining IP links whose restoration paths are longer than their primary path, we plot the top 10 longest restoration paths in Fig. 17(b) (with transponder frequency tuning) and Fig. 17(c) (without transponder frequency tuning). We observe that all restoration fiber paths are shorter than 5,000 km; hence, they can support 100 Gbps modulation based on our device datasheet in Table 6. Higher datarates may be possible for some shorter paths, and restoring such highly-modulated wavelengths to a longer restoration fiber path may trigger the change of modulation formats. Prior work has demonstrated that modulation change latency in WANs is 70 seconds using commodity hardware and can be improved to 35 ms [77].

**WAN transponders datarate vs. reach.** We use the following specifications from our optical device vendors to plan and manage the optical layer at Facebook. For the same wavelength slot, higher capacity is achieved with more aggressive modulation, thus requiring shorter transmission distance [34, 76, 77]. Note that advanced modulation techniques, e.g. probabilistic constellation shaping, could even extend the optical reach with finer-granularity data rates [24].

| Daterate (Gbps) | Reach (km) |
|---|---|
| 100 | 5000 |
| 200 | 3000 |
| 300 | 1500 |
| 400 | 1000 |

**Table 6: Terrestrial long-haul optical transponder specification sheet at Facebook.**

## A.2 Routing and Wavelength Assignment (RWA) Formulation

The RWA is a classical problem in optical networking [88]. Here, we show the RWA formulation used in ARROW.

**Separation of routing and wavelength assignment.** Routing and wavelength assignment problems can be combined into one formulation, but ARROW separates the routing step from the wavelength assignment step. The RWA problem is known to be computational intractable because it requires jointly optimizing a wavelength's routing path length and the wavelength's frequency and modulation assignment. ARROW's separation significantly reduces the formulation complexity, making the problem solvable for large topologies like Facebook within several minutes, while not compromising optimality because routing paths are pre-computed respecting wavelengths' modulation format maximum transmission reach as noted in Table 6. Note that our RWA is solved on a provisioned brown-field (some wavelengths are already populated to carry live traffic) optical network for restoring the failed wavelengths only, while prior proposals are for green-field optical network planning (the network uses dark fibers, and no wavelengths are populated). This is another reason why our RWA can be solved within several minutes.

**Routing the restored wavelengths.** Consider a fiber cut scenario $q$ where a set of IP links $E$ is lost. To restore each IP link $e \in E$ and revive its IP layer capacity, we run $k$-shortest-path algorithm [86] to find its $k$ surrogate fiber paths $\{P_e^1, ..., P_e^k\}$ with path-length upper bounds based on the modulation format of the failed IP links.[10] Note that we allow multiple restored wavelengths of an IP link to be routed over multiple surrogate restoration fiber paths because their IP-layer bandwidth capacity can be aggregated thanks to the link aggregation protocol (IEEE 802.3ad LACP) [44]. We then represent the routing information using a binary parameter $\pi_\phi^{e,k}$ as 1 if a fiber $\phi \in P_e^k$ and 0 otherwise.

**Wavelength assignment of the restored wavelengths.** After obtaining the routing path for restored wavelengths, we need to assign a frequency to each restored wavelength. Consider a fiber $\phi$ on the optical layer topology; its spectrum occupancy can be represented by a binary vector $\phi.spectrum = [0, 0, 1, ..., 1, 0]$ (e.g., 96 wavelength slots under ITU-T DWDM standard [38]), where 0 means this wavelength slot is already utilized by some working wavelengths carrying live traffic, and 1 means this wavelength slot is available for hosting the reconfigured wavelength for optical restoration. Under one failure scenario, we define a binary variable $\xi_{\phi,w}^{e,k}$ as 1 if the restored IP link $e$'s $k$ surrogate restoration fiber path uses wavelength slot $w$ on fiber $\phi$ and 0 otherwise. We further define an integer variable $\lambda_e^k$ that represents the number of restored wavelengths of failed IP link $e$ on its $k$ surrogate restoration fiber paths. Hence, the wavelength assignment of all restored wavelengths for each restored IP link $e$ should follow the following constraints when maximizing the total restored wavelength number $\sum_e \sum_k \lambda_e^k$.

$$\forall \phi, w: \qquad \sum_e \sum_k \xi_{\phi,w}^{e,k} \leq \phi.spectrum[w] \quad (14)$$

$$\forall e, k, \phi: \qquad \lambda_e^k \times \pi_\phi^{e,k} = \sum_w \xi_{\phi,w}^{e,k} \quad (15)$$

$$\forall e, k, w: \qquad \xi_{\phi,w}^{e,k} = \xi_{\phi',w}^{e,k} \quad \textbf{if } \phi, \phi' \in P_e^k \quad (16)$$

$$\forall e: \qquad \sum_k \lambda_e^k \leq \gamma_e \quad (17)$$

Constraint (14) ensures each available wavelength on surrogate restoration fibers can only be used once. Constraint (15) formulates the relationship between restored capacity $\lambda_e^k$ of IP link $e$ routed on its $k$ surrogate restoration fiber path $P_e^k$ and the wavelength assignment on fibers $\phi \in P_e^k$. Constraint (16) is the wavelength

---

[10]Higher-capacity links use more aggressive modulation on the optical layer and hence require shorter transmission distance [34, 76, 77].
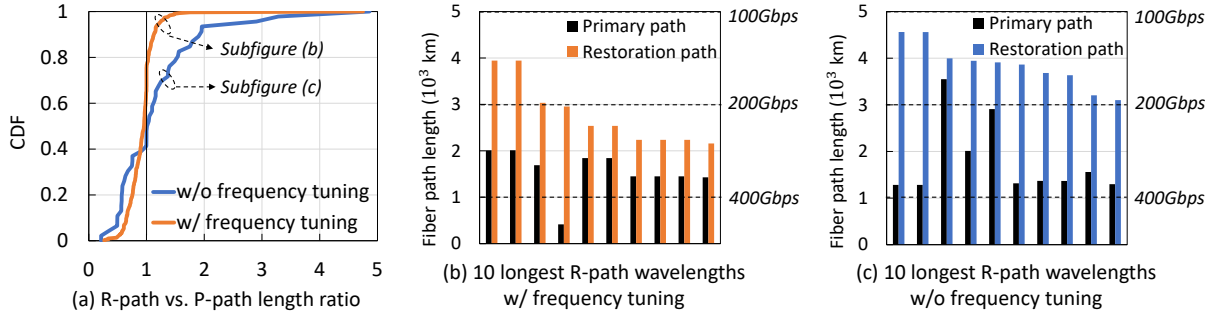
Figure 17: Path length inflation analysis (R-path is short for restoration path, and P-path is short for primary path). 50% of IP links' restoration path is *shorter* than their primary paths. Hence, no modulation change needed. The length ratio between restoration path and primary path has implications for both transponder modulation on the optical layer [50] and service latency on the IP layer [69].

continuity constraint that ensures the occupied wavelength slots are the same all along the fiber path [12]. Constraint (17) ensures failed IP link $e$'s total restored wavelength number on all $k$ surrogate restoration fiber paths $\lambda_e = \sum_k \lambda_e^k$ does not exceed the initial wavelength number $\gamma_e$ of failed IP link $e$.

**Relaxation of the wavelength assignment ILP.** The above wavelength assignment problem has been proven to be NP-hard [11, 12]. To solve this problem in practical network scale for polynomial time, we relax the 0-1 binary variable $\xi_{\phi,w}^{e,k}$ to be a floating-point value between 0 and 1 and transform the wavelength assignment problem into a LP. Hence, the restored capacity $\lambda_e$ for link $e$ also becomes a floating-point number.

**Handling non-fractional $\lambda_e$.** In some cases, the relaxed wavelength assignment LP will return integer numbers for $\lambda_e \in \mathbb{Z}^+$. They do not need rounding processing (desired case in normal rounding problems). However, in our problem setting, where we are generating a set of different restoration candidates, this non-fractional output will result in 0 probability to apply either rounding up or down. Note that in standard randomized rounding techniques, the non-fractional results mean no rounding operation is needed. But in our problem setting, this limits the exploration space of candidate restoration options. To address the non-fractional condition, in our evaluations, when we encounter this situation, we set the probability of rounding up and down as 0.3, and the probability of not rounding as 0.4.

### A.3 Proof of Theorem 3.1

Following the discussion in §3.3, we aim to provide a probabilistic guarantee of LotteryTickets' optimality with randomized rounding. Our proof is based on an the assumption that if the optimal LotteryTicket appears in the set of LotteryTickets as input to Arrow ($z^{opt} \in Z$), Arrow's TE will find the optimal allocation. This assumption is true for Arrow's binary TE formulation (shown in Table 9).

To prove Theorem 3.1, we consider the case where there exists a restoration candidate $z^{opt}$ that maximizes Arrow's TE objective. Arrow's randomized rounding technique (Algorithm 1) returns

a set $Z$ of candidate restoration options (LotteryTickets). Therefore, the probability of $z^{opt} \in Z$ equals the probability of Arrow finding the optimal allocation. We denote this probability as $\rho^q$ in Equation (12). We denote $\kappa$ as the probability of finding $z^{opt}$ with randomized rounding. Hence, for each LotteryTicket, $1 - \kappa$ denotes the probability that this LotteryTicket is not optimal and $(1 - \kappa)^{|Z^q|}$ denotes the probability that all $|Z^q|$ LotteryTickets are not optimal. Therefore, the probability of finding $z^{opt}$ (at least one optimal LotteryTicket) is $1 - (1 - \kappa)^{|Z^q|}$.

To find $\kappa$, we need to derive the probability of finding the best LotteryTicket using our randomized rounding algorithm (Algorithm 1). We apply probabilistic randomized rounding in two steps: 1) rounding stride decision (line 6); 2) rounding up/down decision (line 7). With $\delta$ as the maximum rounding stride, the probability that our chosen stride will be optimal is $1/\delta$. Therefore, for one failed IP link, the probability of obtaining the best LotteryTicket from the initial LP floating point solution is $1/\delta \times \mathbf{Pr}\{round\ up/down\}$. Since each failure scenario $q$ may affect multiple IP links, every LotteryTicket may contain the restoration value of several IP links. Therefore, the probability $\kappa$ of finding the optimal LotteryTicket is determined by multiplying the probabilities of all failed IP links. Hence, we derive Equation (13).

### A.4 Optimal IP/Optical Formulation for Restoration-Aware TE

We present the difference between conventional TE and restoration-aware TE in Fig. 18. As we discussed in §3, the LotteryTicket design enables Arrow to balance computation complexity and solution optimality: 1) TE without optical restoration information (conventional TE that only operate on the IP layer, e.g., FFC, TeaVaR, etc, shown in Fig. 18(a)), and 2) TE with full optical restoration information (Table 7, shown in Fig. 18(b) and Fig. 18(c)). We present the optimal IP/optical version of Arrow formulation in Table 7. In Fig. 18(d), we depict the design of Arrow's approach to avoid the excessive complexity of the joint IP/optical formulation based on abstracting optical layer with LotteryTickets.
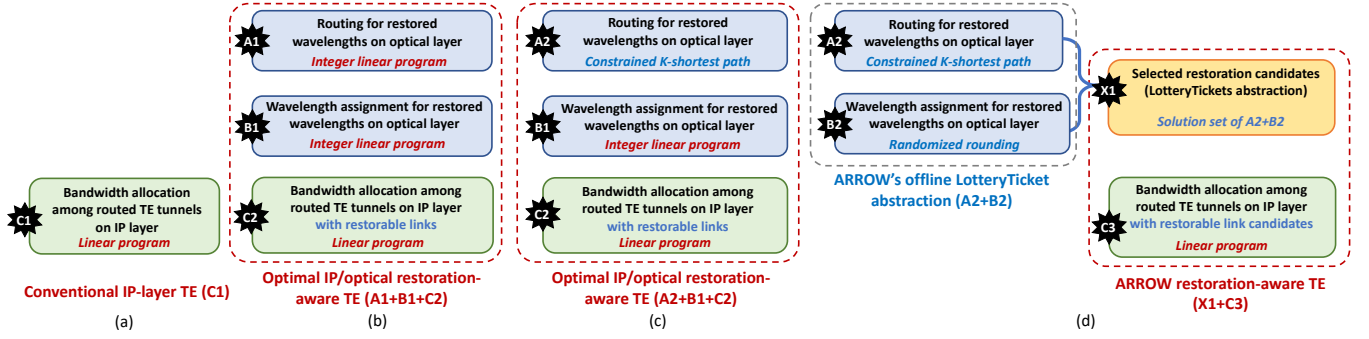
**Figure 18: Illustrative relationship of conventional TE formulation(subfigure a), optimal IP/optical restoration-aware TE formulation (subfigures b, c), and ARROW TE formulation (subfigure d).**

**IP/optical TE is optimal.** Ideally, an IP/optical cross-layer TE takes a set of failure scenarios as well as full optical-layer information for optical restoration into account while optimizing traffic allocations on tunnels to maximize overall network throughput. This IP/optical cross-layer TE is optimal in terms of network throughput, because IP-layer flow allocation is jointly optimized with optical-layer wavelength restoration.

**Optimal IP/optical TE is computationally intractable.** As we shown in Fig. 18(b), the cross-layer TE problem takes a set of failure scenarios as input, and for each failure scenario it also contains constraints on the routing (A1) and wavelength assignment (B1) problems on the optical layer following all the optical-layer constraints (e.g., wavelength continuity). This is a RWA problem, which has been proven to be NP-Complete [88]. Moreover, as shown in Fig. 18(c), another way to formulate the cross-layer TE is to convert flow conservation constraints of the wavelength routing problem into a set of pre-calculated surrogate restoration fiber paths (A2) and feed them into the optimization for wavelength assignment (B1). This operation reduces the problem size, but the selection of which surrogate restoration fiber path to use and which wavelength to restore is still an integer problem.

**Dynamic restorable tunnels.** Unlike the ARROW formulation in Table 2 and Table 3, the state that whether a tunnel is restorable or not is not an input to the optimization, but being dynamically decided internally during the optimization. This is because which IP link can be restored on the optical layer is jointly decided with TE flow allocation, not as input (e.g., LotteryTickets). Therefore, in the joint IP/optical formulation, we do not have a restorable tunnels set $Y_f^{z,q}$ as input. Instead, we take all failed tunnels $T_f - T_f^q$ of flow $f$ as candidates to be potentially restored and jointly make the decision with the optical-layer RWA.

**Size of the optimal formulation.** The optimal formulation is computational intractable, and cannot yield solutions within reasonable amount of time. Therefore, we list the size of the optimization formulation in Table 8. As we can find in Table 8, for an ILP with such problem size, it is not feasible to obtain an optimal solution even in several days with state-of-the-art optimization solvers, e.g., Gurobi [66]. Hence, the cross-layer optimal formulation is computationally intractable for TE.

**Constraints.** Constraints (18-20) are identical as Constraints (1-3) in Table 2 and Constraints (7-9) in Table 3. Constraint (21) considers all failed tunnels as potential dynamic restorable tunnels and ensure that the total allocated bandwidth for flow $f$ should be no larger than sum bandwidth of its residual tunnels and dynamic restorable tunnels. Constraint (22) ensures the total bandwidth capacity of restorable tunnels does not exceed the restorable bandwidth of the failed link $e$ under failure scenario $q$. Constraints (23-26) are similar as Constraints (14-17) and ensure RWA constraints for restoring wavelengths of each failed IP link under failure scenario $q$. Constraint (27) connects the total number of restorable wavelength $\sum_k \lambda_{e,q}^k$ of link $e$ with the bandwidth capacity $r_e^q$ of restorable IP links $e$ under failure scenario $q$ by multiplying each wavelength with its modulation format.

## A.5 Binary ILP Formulation for ARROW's TE with LotteryTickets

ARROW's ticket selection process can be formulated as a binary-integer linear problem due to ticket selection (represented by a binary variable $x^{z,q}$), as shown in Table 9. The advantage of this binary formulation is that it can confirm the assumption in Proof 3.1 (if the optimal LotteryTicket appears in the set of LotteryTickets as input to ARROW ($z^* \in Z$), ARROW's TE finds the optimal allocation), however, at the cost of computational complexity.

**Constraints.** Constraints (28-30) are identical to Constraints (1-3) in Table 2, Constraints (7-9) in Table 3 and Constraints (18-20) in Table 7. Constraint (31-32) are augmented from Constraints (4-5) with the binary selection of LotteryTicket. Constraint (33) ensures only one LotteryTicket can be selected for each failure scenario $q$.

## A.6 ROADM Reconfigurations

**Parallel ROADM configurations.** In practice, the surrogate restoration fiber path includes multiple ROADMs, and all of them need to be configured for data/noise replacement. In ARROW, we avoid serially configuring each of these ROADMs. Instead, we group all ROADMs into two categories: *add/drop* ROADMs representing source/destination sites and *intermediate* ROADMs representing ROADMs that act as optical switches to steer light to a designated direction. ARROW reconfigures all ROADMs in each group in parallel: it first reconfigures all add/drop ROADMs as well as their

| | | |
|---|---|---|
| **Standard Optical-Layer Input Parameters** | $G(\Psi, \Phi)$ | Optical-layer network graph with ROADM set $\Psi$ and fiber set $\Phi$. |
| | $\phi.spectrum[w]$ | Each $\phi \in \Phi$ contains a binary vector indicating if wavelength slot $w$ of this fiber $\phi$ is occupied or not. |
| | $Q = \{q\}$ | Fiber cut failure scenarios. Each $q \in Q$ is represented in a $|\Phi|$-size binary vector $\{..., h_\phi^q, ...\}, \phi \in \Phi$ indicating the healthy state of each fiber. Using the provisioned mapping between IP links and optical fibers, we can derive another $E$-size binary vector $\{..., h_e^q, ...\}, e \in E$ indicating whether each IP link is affected. |
| | $\pi_\phi^{e,k}$ | A binary parameter, 1 if IP link $e$'s $k$-th optical-layer surrogate restoration fiber path traverses fiber $\phi$, otherwise 0. |
| | $\gamma_e$ | Number of wavelengths of IP link $e$ before failure. |
| | $P_e^k$ | Failed IP link $e$'s $k$ surrogate restoration fiber path for restoration. |
| **TE Input** | Table 1 | Standard TE input parameters. |
| **Optical Output** | $\xi_{\phi,w}^{e,k,q}$ | Binary variable, if IP link $e$'s $k$ surrogate restoration fiber path is routed on fiber $\phi$ using wavelength $w$ under scenario $q$. |
| | $\lambda_e^{k,q}$ | Integer variable, number of restored wavelengths on $k$ surrogate restoration fiber path of IP link $e$ under scenario $q$. |
| | $r_e^q$ | Restorable bandwidth capacity for link $e$ under scenario $q$. |
| **TE Output** | $b_f$ | Total allocated bandwidth for flow $f$. |
| | $a_{f,t}$ | For flow $f$, the allocated bandwidth on tunnel $t \in T_f$. |

**Maximize:** $\sum_{f \in F} b_f$
**Subject to:**

$$\forall f: \quad \sum_{t \in T_f} a_{f,t} \geq b_f \tag{18}$$
$$\forall e: \quad \sum_{f \in F} \sum_{t \in T_f} a_{f,t} \times L[t,e] \leq c_e \tag{19}$$
$$\forall f: \quad 0 \leq b_f \leq d_f \tag{20}$$
$$\forall f,q,: \quad \sum_{t \in T_f - T_f^q} a_{f,t} + \sum_{t \in T_f^q} a_{f,t} \geq b_f \tag{21}$$
$$\forall e,q: \quad \sum_{f \in F} \sum_{t \in Y_f^{z,q}} a_{f,t} \times L[t,e] \leq r_e^q \tag{22}$$
$$\forall \phi, w, q: \quad \sum_e \sum_k \xi_{\phi,w}^{e,k,q} \leq \phi.spectrum[w] \tag{23}$$
$$\forall e,k,\phi,q: \quad \lambda_{e,q}^k \times \pi_\phi^{e,k} = \sum_w \xi_{\phi,w}^{e,k,q} \times h_\phi^q \tag{24}$$
$$\forall e,k,w,q: \quad \xi_{\phi,w}^{e,k,q} = \xi_{\phi',w}^{e,k,q} \quad \text{if } \phi, \phi' \in P_e^k \tag{25}$$
$$\forall e,q: \quad \gamma_e \times h_e^q \leq \sum_k \lambda_e^{k,q} \leq \gamma_e \tag{26}$$
$$\forall e,q: \quad \sum_k \lambda_{e,q}^k \times \lambda_{e,q}^k.modulation = r_e^q \tag{27}$$

**Table 7: Arrow joint IP/optical TE formulation.**

noise sources, and then switches to reconfiguring all intermediate ROADMs and noise sources.

**Number of ROADMs to be reconfigured.** Reconfiguring wavelengths from the cut fiber path to the surrogate restoration fiber

| Topology | # of binary vars. | # of continuous vars. | # of constraints |
|---|---|---|---|
| Facebook | 12,280 million | 72 thousand | memory overflow |
| IBM | 81 million | 6.5 thousand | 192 million |
| B4 | 52 million | 3.5 thousand | 119 million |

**Table 8: Size of joint IP/optical TE formulation.**

| | | |
|---|---|---|
| **Arrow Binary ILP Input Parameters** | Table 1 | Standard TE input parameters. |
| | $Z^q = \{z\}$ | Set of LotteryTicket indexes under scenario $q$. |
| | $r_e^{z,q}$ | Restorable bandwidth capacity for link $e$ under scenario $q$ and LotteryTicket $z$. |
| | $Y_f^{z,q}$ | Restorable tunnels for flow $f$ under scenario $q$ and LotteryTicket $z$. |
| | $M$ | A big number. |
| **Arrow Binary ILP Output** | $b_f$ | Total allocated bandwidth for flow $f$. |
| | $a_{f,t}$ | For flow $f$, the allocated bandwidth on tunnel $t \in T_f$. |
| | $x^{z,q}$ | 1 if the $z$ LotteryTicket is selected under scenario $q$, otherwise 0. |

**Maximize:** $\sum_{f \in F} b_f$
**Subject to:**

$$\forall f: \quad \sum_{t \in T_f} a_{f,t} \geq b_f \tag{28}$$
$$\forall e: \quad \sum_{f \in F} \sum_{t \in T_f} a_{f,t} \times L[t,e] \leq c_e \tag{29}$$
$$\forall f: \quad 0 \leq b_f \leq d_f \tag{30}$$
$$\forall f,q,z: \quad \sum_{t \in Y_f^{z,q}} a_{f,t} + \sum_{t \in T_f^q} a_{f,t} \geq b_f - M(1 - x^{z,q}) \tag{31}$$
$$\forall e,q,z: \quad \sum_{f \in F} \sum_{t \in Y_f^{z,q}} a_{f,t} \times L[t,e] \leq r_e^{z,q} + M(1 - x^{z,q}) \tag{32}$$
$$\forall q: \quad \sum_{z \in Z^q} x^{z,q} = 1 \tag{33}$$

**Table 9: Arrow TE binary ILP formulation.**

path requires the ROADMs and ASE noise sources (at both add-/drop nodes and intermediate nodes) to be reconfigured. We quantify the number of devices to be reconfigured with Arrow for every fiber on Facebook's optical backbone. Fig. 19 shows the CDF of number of add/drop and intermediate ROADMs to be reconfigured withunder Arrow at Facebook. We observe that for 80% of the fiber cut events, the number of add/drop ROADMs is less than 10, while the number for intermediate ROADMs is less than 6. The reason why there may be more than 2 *Add/drop* ROADMs is that the failed wavelengths on a cut fiber do not necessarily originate and terminate at the endpoints of the broken fiber. Their source and destination sites could be any other ROADM sites on the optical layer.

## A.7 Wavelength Reconfiguration in Legacy Optical Layer is Slow

**Wavelength reconfiguration is non-trivial.** The presence of amplifiers on fibers introduces a non-trivial challenge to wavelength reconfiguration because of the complex relationship between amplifiers' gain control mechanism and the wavelengths traversing the fiber. A sudden change to the set of wavelengths, i.e., adding or removing wavelengths simultaneously, could result in unpredictable power fluctuations on each amplifier which, in turn, could lead to packet loss/errors in the IP layer. Current device manufacturers and backbone operators have settled on a conservative stabilization
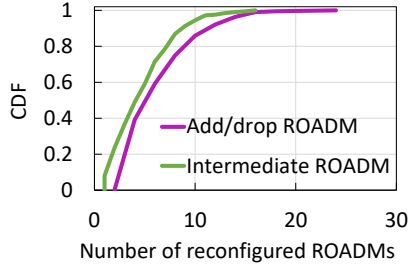
**Figure 19: Number of ROADMs that need to be reconfigured for each fiber cut in Facebook's WAN.**
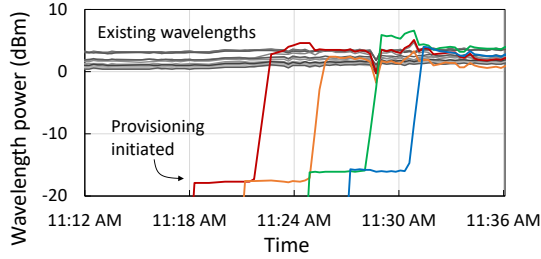


**Figure 20: Amplifiers take minutes to adjust power during wavelength reconfiguration on a 2,000 km fiber path with 24 cascaded amplifier sites between Canada and US.**
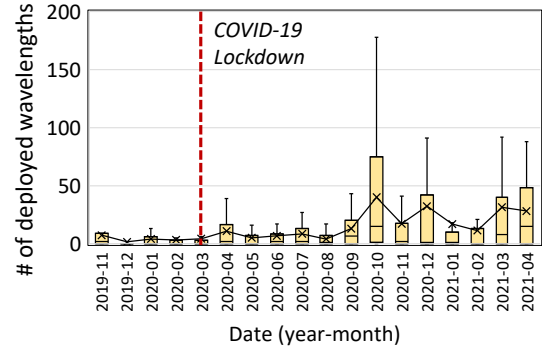


**Figure 21: Monthly wavelength deployment at Facebook.**



**Figure 22: (a) CDF of number of IP links per fiber. (b) CDF of number of wavelengths per IP link.**

process with multiple *observe-analyze-act* control loops that takes several minutes per amplifier. To understand the current practice, we shadowed the process of reconfiguring four wavelengths on a 2,000 km fiber path with 24 cascaded amplifier sites between Canada and US in Facebook. This reconfiguration was part of an automatic maintenance procedure. As shown in Fig. 20, it takes 14 minutes to reconfigure four wavelengths using *legacy hardware* without noise loading at Facebook.

**Slow wavelength reconfiguration is common in production.** Even though wavelength reconfiguration is a slow process, it is an essential part of a production backbone's operations. Fig. 21 shows the number of monthly deployed wavelengths from November 2019 to April 2021 in Facebook. We observe that during COVID-19, more wavelengths were deployed since March 2020 to handle the increase of online traffic at Facebook [85]. Arrow's fast wavelength reconfiguration can bring advantages to wavelength deployments as well.

### A.8 IP-to-Optical Topology Mapping

As discussed in Section 6, in large-scale WANs, the IP-layer topology tends to be denser than the optical-layer topology [65, 71]. Fig. 22 shows the CDF of number of IP links per fiber and number of wavelengths per IP link in Facebook. We use these distributions to guide us to generate the IP-layer topologies for B4, IBM, and Facebook.

### A.9 Prior Work Comparison

Three main techniques are used to mitigate the impact of lost capacity caused by fiber cuts: (*i*) failure-aware TE; (*ii*) optical path protection; (*iii*) optical restoration. Table 10 illustrates the key differences between these approaches and Arrow.

In TE-based solutions (first row in Table 10), Port1 and Port2, as well as Transponder1 and Transponder2, need to be pre-allocated in the WAN to enable the TE to quickly switch traffic from the cut fiber to the failover path. Hence, once a fiber cut occurs, the router ports and transponders associated with the failed fiber become unusable and sit idle until the fiber is repaired.

In OTN-based solutions (the second row in Table 10), Port1 remains active during the fiber cut repair from the IP layer's perspective, but Transponder1 will become idle as traffic is shifted to Transponder2. The OTN approach can save on router ports, but it still requires an extra optical transponder to be in working state so as to be prepared for fast failover under failures.

The third row of Table 10 presents optical restoration [30, 32, 45, 49, 52, 56, 72, 74, 82]. Instead of pre-allocating failover paths, optical restoration techniques dynamically shift the wavelengths from the cut fiber onto healthy *surrogate* fibers. These approaches do not require pre-allocating router ports or transponders. Instead, they leverage optical devices, such as ROADMs, to shift the wavelengths after the fiber cut. The core idea of the optical restoration technique was proposed two decades ago [30] and it remains a popular solution in the optics community as it does not leave router ports or transponders idle during fiber cuts.

Arrow solves two challenges that make previously proposed optical restoration techniques inefficient in large-scale WANs. First, prior works do not consider the interplay between partially restorable IP links and failure-aware traffic engineering. As a result, they may choose sub-optimal restoration candidates. To solve this challenge

| Approach | Failover plan configuration | Failover latency | Practical | Illustration |
|---|---|---|---|---|
| Failure-aware Traffic Engineering [17, 48, 58, 63] | Routing table | $O(ms)$ | ✓ | Port1 and Trans.1 are idle during fiber cut repair |
| Optical Path Protection [14, 19, 59, 68, 80, 81, 84] | OTN configuration | $O(ms)$ | ✓ | Trans.1 is idle during fiber cut repair |
| Optical Restoration [30, 32, 45, 52, 56, 72, 74, 82] | ROADM configuration | 10s mins | ✗ | No idle resources (Wavelengths are dynamically shifted during fiber cut repair) |
| Arrow (this paper) | Routing table, ROADM configuration | $O(s)$ now, $O(ms)$ future | ✓ | No idle resources (Wavelengths are dynamically shifted during fiber cut repair) |

Table 10: Comparison with prior work.

Arrow augments today's TE formulations to consider multiple partial restoration candidates for optimizing IP-layer network throughput (§3). Second, the failover latency of prior proposals is tens of minutes because they require amplifiers' gains to be adjusted. Arrow uses a *noise source* to fully populate the amplifiers' spectrum to bypass the reconfiguration time, achieving an end-to-end failover latency of eight seconds on a WAN-scale testbed (§5). The last row of Table 10 illustrates Arrow's differences from prior approaches.

## A.10 Extensions

**Supporting next-generation C+L optical systems.** There is a trend to expand the C-band spectrum in the optical layer to L-band to scale the network capacity [75]. because of the efficient

abstraction of LotteryTickets, Arrow's TE is orthogonal to optical transmission techniques. The noise loading technique in Arrow's optical layer can smoothly support the expansion of the L band by loading it with noise [83]. As recent advancements of L-band tunable transponders [26] and reconfigurable WSS [36] are becoming commercially available, Arrow can be easily extended to the L-band and support future WANs.