

Scheduling Techniques for Hybrid Circuit/Packet Networks

He Liu^{§*}, Matthew K. Mukerjee[†], Conglong Li[†], Nicolas Feltman[†],
George Papan, Stefan Savage, Srinivasan Seshan[†], Geoffrey M. Voelker,
David G. Andersen[†], Michael Kaminsky[‡], George Porter, and Alex C. Snoeren

University of California, San Diego [†]Carnegie Mellon University [‡]Intel Labs [§]Google, Inc.

ABSTRACT

A range of new datacenter switch designs combine wireless or optical circuit technologies with electrical packet switching to deliver higher performance at lower cost than traditional packet-switched networks. These “hybrid” networks schedule large traffic demands via a high-rate circuits and remaining traffic with a lower-rate, traditional packet-switches. Achieving high utilization requires an efficient scheduling algorithm that can compute proper circuit configurations and balance traffic across the switches. Recent proposals, however, provide no such algorithm and rely on an omniscient oracle to compute optimal switch configurations.

Finding the right balance of circuit and packet switch use is difficult: circuits must be reconfigured to serve different demands, incurring non-trivial switching delay, while the packet switch is bandwidth constrained. Adapting existing crossbar scheduling algorithms proves challenging with these constraints. In this paper, we formalize the hybrid switching problem, explore the design space of scheduling algorithms, and provide insight on using such algorithms in practice. We propose a heuristic-based algorithm, Solstice that provides a $2.9\times$ increase in circuit utilization over traditional scheduling algorithms, while being within 14% of optimal, at scale.

CCS Concepts

•Networks → Bridges and switches; Packet scheduling; Data center networks;

Keywords

circuit networks; packet networks; hybrid networks

*Work done while at UCSD

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

CoNEXT '15, December 01–04, 2015, Heidelberg, Germany

© 2015 Copyright held by the owner/author(s).

ACM ISBN 978-1-4503-3412-9/15/12.

DOI: <http://dx.doi.org/10.1145/2716281.2836126>

1. INTRODUCTION

Today’s datacenters aggregate tremendous amounts of compute and storage capacity, driving demand for network switches with ever-increasing port counts and line speeds. However, supporting these demands with existing packet switching technology is becoming increasingly expensive—in cost, heat, power, and cabling. Packet switches are flexible, capable of making forwarding decisions at the granularity of individual packets. In common modern scenarios, however, this flexibility is unnecessary: many (often consecutive) packets are sent to the same output port. Two key factors contribute to this traffic pattern. First, traffic inside a datacenter often has high spatial locality, where a large fraction of the traffic that enters each switch port is destined for only a small number of output ports [16,23]. Second, traffic is often bursty, with significant temporal locality between packets sharing the same destination [17,23]. The consequence of these two factors is that the traffic demand matrix at a datacenter switch is often both skewed and sparse [5,13,15].

Researchers have seized upon these observations to propose hybrid datacenter network architectures that offer higher throughput at lower cost by combining high-speed optical [4,6,28] or wireless [13,15,31] circuit switching technologies with traditional electronic packet switches. Typically, the circuit switch has a significantly higher data rate than the packet switch, but incurs a non-trivial reconfiguration penalty. While the potential cost savings that hybrid techniques could realize is large, the design space of scheduling algorithms that enable high utilization in hybrid networks is not yet well understood. Earlier work that considers circuit switches with substantial reconfiguration delay offers no guidance about how to negotiate the trade-off between remaining in the current (potentially sub-optimal) circuit configuration vs. incurring a costly reconfiguration delay to switch to a potentially better circuit configuration [6,27,28]. The reconfiguration cost of these systems was so high that they were forced to keep a configuration pinned up for a relatively long period anyway.

In recent years, however, the switching time of optical circuit switches has improved substantially [22]. As a result, an efficient scheduling algorithm for a modern hybrid design must determine: 1) a set of circuit configurations (which ports are connected to which other ports and how

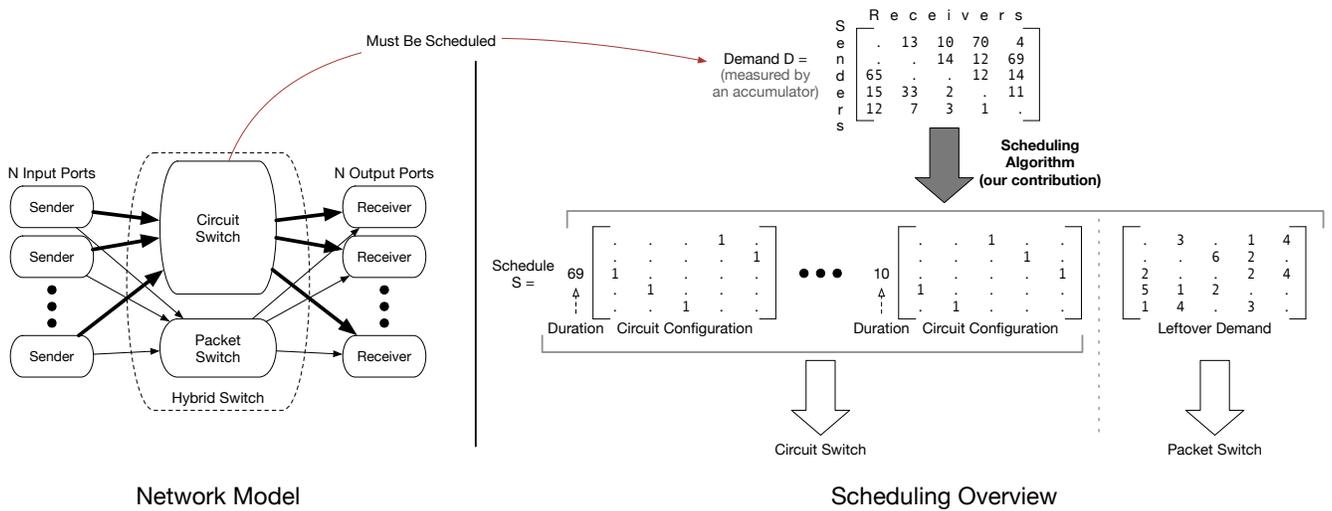


Figure 1: Our model of a hybrid switch architecture and the scheduling process. The circuit switch has high bandwidth, but slow reconfiguration time. The packet switch has low bandwidth (e.g., an order of magnitude lower), but can make forwarding decisions per-packet.

long that configuration should remain in effect) designed to maximize the traffic serviced over the high-bandwidth but slow-to-reconfigure circuit-switched network, and 2) what traffic should be sent to the low-bandwidth but flexible packet switch.

Computing an optimal set of circuit configurations to maximize circuit-switched utilization has no known polynomial-time algorithms, scaling as $O(n!)$ in the number of switch ports (§3). The challenge arises due to the non-trivial switching time between configurations, which necessitates not only sending as much traffic as possible, but doing so in the fewest number of configurations.

The end goal of this paper is an effective and fast heuristic algorithm that delivers high switch utilization. To this end, we first provide a detailed characterization of the problem plus an optimal (but impractical) solution that sheds light on how to design an effective heuristic. We then present our heuristic, Solstice, which provides $2.9\times$ higher utilization compared to previous algorithms by taking advantage of the known sparsity and skew of datacenter workloads—some of the same features that make the traditional scheduling problem hard.

The contributions of this work are as follows:

1. Characterizing the hybrid switch scheduling problem.
2. Exploring the design space of hybrid scheduling:
 1. **Lower bound:** an instantly computable but loose bound on the minimum amount of time it takes to serve all demand (but provides no actual schedules).
 2. **Optimal scheduling:** optimally schedule all demand with minimal time; impossible to run in real time at scale.
 3. **Heuristic algorithm (“Solstice”):** runs in real time at scale, but slightly underperforms optimal (by at most 14% at target scale).

4. **Heuristic + optimization (“Solstice++”):** runs at scale (though not in real time), but tightens the gap between Solstice and optimal (at most 12% from optimal at target scale).

3. Insight into the challenges and benefits of using hybrid switches, with a focus on high circuit utilization.

2. BACKGROUND

We consider a single switch in a hybrid network fabric that consists of n ports. In the context of datacenters, these ports would typically connect to individual servers or Top-of-Rack (ToR) switches. We leave multiple-switch networks to future work. Our model assumes each port is logically input queued. In some realizations, the queues are located at the senders themselves [19], although alternatively, the queues could be located at the ToR switches or the hybrid switch itself.

Our abstraction of a hybrid switch (shown on the left-hand side of Figure 1) consists of two separate switches: a circuit switch, typically optical or RF, capable of forwarding at very high bandwidth, and a low-bandwidth (e.g., an order of magnitude lower) packet switch. Both switches source packets from the queues at each of the n input ports, structured as virtual output queues (VOQs). Although the circuit switch has a significantly faster data rate than the packet switch, it incurs a non-trivial reconfiguration penalty.

Prior work has focused on building such a switch [6, 19, 22, 28], with little focus on how to schedule traffic, instead relying on an omniscient oracle to compute optimal switch configurations. ReacToR, for example, leaves the selection and evaluation of a hybrid scheduler as future work [19].

2.1 Switch model

In our model, each input port of the hybrid switch is simultaneously connected to both the packet switch and the circuit switch. At any point in time, however, at most one VOQ at each input port may be serviced by the circuit switch, whereas

Symbol	Definition
n	number of switch ports
δ	circuit reconfiguration time
r_c, r_p	circuit/packet link rates
D	input demand matrix ($n \times n$)
E	demand sent to packet switch ($n \times n$)
P_i	circuit switch configuration ($n \times n$)
t_i	time duration of P_i
m	number of configurations

Table 1: Notation used throughout the paper.

multiple VOQs may be drained simultaneously by the packet switch. The circuit switch functions as a crossbar: it can connect any input port to any output port, but no output port may be connected to multiple input ports, and no input port may be connected to multiple output ports (aside from their connection to the packet switch) in a single configuration.

The circuit switch can be reconfigured with the cost of a fixed time delay δ (e.g., 20 μ s [19]). Some technologies allow circuits that do not change during a reconfiguration to forward data during the reconfiguration period. We assume a pessimistic view that all communication stops during a reconfiguration, allowing our scheduler to function with a wider set of technologies. The packet switch, on the other hand, can service traffic at all times.

To ensure high circuit utilization, each circuit configuration must remain active for a long period with respect to δ . For example, to ensure 90% link utilization over the circuit switch, the average duration of a configuration needs to be at least 9 δ (e.g., 180 μ s) to amortize the reconfiguration delay.

One important distinction between our model and traditional switches is that there is no queueing at the output ports of the circuit switch. This restriction rules out any crossbar scheduler that requires a speed-up factor. Hybrid switches instead use a lower-data-rate commodity packet switch (without constraints on queueing/speed-up), to make up for the reconfiguration delay and any scheduling inefficiency. We will see that this addition provides an improvement compared to existing approaches.

2.2 Formalizing the problem

Our goal is to calculate a schedule for the circuit switch, and to determine what data to send to the packet switch, such that we service all demand (i.e., no starvation) while maximizing utilization. How the scheduler learns about the traffic demand is orthogonal to this work, but some possibilities include estimation/prediction algorithms or simply accumulating the demand before transmission.

2.2.1 Notation

In order to formalize our goal, we define some of the core concepts. We summarize our notation in Table 1 and below:

The hybrid switch: Our hybrid switch has n full-duplex ports. The circuit switch has a reconfiguration time of δ seconds and a link rate of r_c bits/second. The packet switch has a link rate of $r_p \ll r_c$ (e.g., 1:10) bits/second.

Formula	Definition
Controllable Variables:	
m	number of configurations
P_i	circuit switch configuration
t_i	time duration of P_i
E	demand sent to packet switch
Goal:	
$\min T = (\sum_{i=1}^m t_i) + m\delta$	minimize total time
Constraints:	
1) $E + \sum_{i=1}^m r_c t_i P_i \geq D$	demand satisfaction
2) $\forall i : \sum_{j=1}^m E_{i,j} \leq r_p T$	cap outbound packet links
3) $\forall j : \sum_{i=1}^m E_{i,j} \leq r_p T$	cap inbound packet links

Table 2: Summary of the problem.

Demand: We express demand as a matrix D of size $n \times n$, where the rows are sources and the columns are destinations. $D_{a,b} \in \mathbb{R}_0^+$ is the amount of data port a wants to send to port b , in bits. In the resulting schedule, some of this demand, which we denote by matrix E , will be sent to the packet switch; E is an $n \times n$ matrix, where $E_{a,b}$ is the portion of the demand $D_{a,b}$ sent from a to b via the packet switch, in bits.

Scheduling: A circuit switch schedule is a set of configurations $\{P_1, P_2, \dots, P_m\}$ and an associated set of durations $\{t_1, t_2, \dots, t_m\}$. Each configuration P_i is an $n \times n$ binary matrix encoding which nodes are connected to each other in the circuit switch. $P_{i_{a,b}} = 1$ iff port a can send to port b during this configuration. Because the circuit switch connects each sender to exactly one receiver and vice-versa, all P_i are permutation matrices (i.e., have exactly one 1 in each row/column). Each configuration P_i is associated with a corresponding duration t_i that indicates how long the circuit switch should remain in that configuration.

2.2.2 Overall goal

Our goal is to minimize the amount of transmission time it takes to schedule all demand, thus maximizing utilization. We wish to fully schedule all demand before considering new demand to ensure fairness and to avoid starvation. To achieve this, our algorithm selects the circuit switch schedule (m, P_i, t_i) as well as which data to forward via the packet switch (E). This process is depicted in the right-hand side of Figure 1. Further, we need to formally define our goal, *total time*, and two constraints, *demand satisfaction* and *packet switch capacity*, which we do in Table 2 and below:

We define *total time* as the amount of time scheduled on the circuit switch plus the amount of time spent switching configurations:

$$T = \left(\sum_{i=1}^m t_i \right) + m\delta.$$

Time used by the packet switch is constrained to be concurrent with the time spent on the circuit switch, limiting the

packet switch capacity. Thus, for each outbound link i , the amount of data admissible is:

$$\sum_{j=1}^m E_{i,j} \leq r_p T.$$

Inbound links are constrained similarly.

Demand is satisfied when, for each source/destination pair, the amount of data served on the packet switch plus the amount served on the circuit switch is greater than or equal to the demand:

$$E + \sum_{i=1}^m r_c t_i P_i \geq D.$$

2.3 Demand matrices

A key assumption in our work is that demand matrices are sparse and skewed. We now discuss both assumptions.

Sparsity: For a demand matrix D , counting the number of non-zero elements in each row and column results in $2n$ values. The largest of these values we refer to as D_{count} . D is “sparse” when D_{count} is small. Sparse matrices can be scheduled more efficiently on a circuit switch since they inherently require fewer configurations. In practice, for a fixed period of demand accumulation, D_{count} has been shown to be bounded by a constant (≈ 5) in an empirical study [16]. More recent work has suggested that D_{count} has grown larger (e.g., low 10s [23]), but it appears that D_{count} is growing much slower than n .

Skewness: A matrix is “skewed” if the ratio between its maximum and minimum non-zero elements of the matrix is high. Assuming a fixed period of demand accumulation, skewed demand matrices are fundamentally less efficient for a circuit-switched network to serve, since as the magnitude of small demands decreases, the durations of the circuit configurations required to service those demands becomes short relative to the reconfiguration time, decreasing overall utilization. In contrast, for hybrid networks, very small elements in the demand are likely well served by the packet—rather than circuit—switch.

3. OPTIMALITY

To better understand our heuristic algorithm’s results, we construct an integer linear program (ILP) that computes an optimal schedule. Though it is impractical for online use or at scale, it effectively considers all possible permutation matrices to determine which to use and for how long and it provides a useful and exact lower bound for comparison with other approaches.

3.1 Formulation

Starting from a large candidate set of circuit switch configurations (i.e., all $n!$ binary permutation matrices), $\{P_1, \dots, P_{n!}\}$, our goal is to compute $\{t_1, \dots, t_{n!}\}$, the time spent in each (potential) configuration P_i . Note that with a large candidate set of configurations and a sparse demand matrix, almost all t_i will likely be zero, meaning that the corresponding P_i are not used by the resulting schedule. We

$$\min \left(\sum_{i=1}^{n!} t_i \right) + m\delta$$

subject to:

- 1) $E + \sum_{i=1}^{n!} r_c t_i P_i \geq D$
- 2) $\forall i : \sum_{j=1}^n E_{i,j} \leq r_p T$
- 3) $\forall j : \sum_{i=1}^n E_{i,j} \leq r_p T$
- 4) $\forall i : t_i \leq \max(D) l_i$

Figure 2: An ILP to find optimal schedules for a hybrid switch.

define a binary indicator variable l_i that denotes whether configuration P_i is employed by the solution (i.e., its corresponding t_i is non-zero). The number of configurations used in the schedule, m , is then

$$m = \sum_{i=1}^{n!} l_i.$$

The remainder of the demand matrix D not serviced by the m selected circuit switch configurations forms E , the $n \times n$ matrix served by the packet switch.

Figure 2 shows the ILP, which minimizes the total length of the schedule (i.e., duration of the configurations plus the switching overhead; T in Table 2) subject to four constraints. The first three are effectively identical to those in Table 2. The final technical constraint ensures we incur a reconfiguration delay only for permutations included in our final schedule.

3.2 Candidate permutations

An obvious challenge with this approach is that it considers all $n!$ possible circuit configurations. Although modern ILP solvers (e.g., Gurobi [11]) are very fast, $n!$ is impractical for n on the order of modern switch port counts (e.g., at least 48). Even considering all possible configurations for a 16-port switch would require more than 4 petabytes of memory.

Fortunately, it is possible to consider only a much smaller set of configurations and still maintain optimality. We observe that for a given (sparse) demand matrix D , most possible circuit configurations connect two nodes with no demand. Removing these “useless” links yields a partial configuration we refer to as a *class*. Many configurations yield the same class and, thus, are redundant; we need to keep only one example from each class. Moreover, when comparing two classes, one class may be a strict superset of another, meaning the subset class is redundant, so we can remove it as well.

We employ a straightforward dynamic programming-based algorithm (omitted for space) to generate an example configuration from each class. Although we avoid generating all $n!$ permutation matrices, we find that class generation still produces $O(n!)$ candidate permutations—though it does so with a constant speedup of $\sim 100\times$. Even so, solving a small-scale (e.g., 12-port) ILP still takes around 5 minutes on our hardware (see Table 3 in §5), and thus we require an approximation algorithm that can provide nearly optimal results quickly at scale.

Algorithm 1: Birkoff-von Neumann decomposition

```

input :  $k$ -bistochastic matrix  $D$ 
         link rate for circuit switch:  $r_c$ 
output:  $m$  circuit configurations and durations:  $\{P_i\}, \{t_i\}$ 
 $i \leftarrow 1$ 
while  $D > 0$  do
   $B \leftarrow \text{BinaryMatrixOf}(D)$ 
  Interpret  $B$  as a bipartite graph of senders to receivers.
  Find a perfect matching  $P_i$  of  $B$ .
  Interpret  $P_i$  as a permutation matrix.
   $t_i \leftarrow \min\{D_{a,b} \mid P_{i,a,b} = 1\}/r_c$ 
   $D \leftarrow D - r_c t_i P_i$ 
   $i \leftarrow i + 1$ 
end
 $m \leftarrow i - 1$ 

```

4. SOLSTICE

The classical approach to scheduling crossbar circuit switches is Birkoff-von Neumann decomposition (BvN) [3]. BvN produces a schedule of at most n^2 configurations and associated durations that fully satisfies the given demand. BvN, however, is inappropriate for our target network environment for two reasons. First, BvN uses solely the circuit switch, whereas a hybrid scheduler must effectively use the packet switch. Second, BvN does not try to minimize the number of configurations, which is necessary in practice to avoid expensive reconfiguration delays. We address these shortcomings by developing Solstice, a heuristic-based scheduling algorithm targeted for the hybrid case. Solstice is closely related to BvN, so we first review BvN scheduling.

4.1 Birkoff-von Neumann scheduling

The Birkoff-von Neumann theorem forms the theoretical underpinning of BvN decomposition [3]. BvN assumes a non-negative square input matrix D that is k -bistochastic, meaning that each row and column of D sums to the same value k . The BvN theorem states that all k -bistochastic matrices can be decomposed into a set of at most n^2 permutation matrices and corresponding non-negative durations which sum to k . The steps for doing so are shown in Algorithm 1. While real demand matrices are not typically k -bistochastic, a pre-processing step can make them k -bistochastic (e.g., Sinkhorn’s algorithm [24]) by adding artificial demand.

One alternative to inserting artificial traffic into the demand matrix would be to employ an algorithm that could decompose non- k -bistochastic matrices by considering partial circuit configurations (i.e., configuring only a subset of ports). Such an algorithm would have to consider all possible subsets of the n^2 port pairs (of which there are 2^{n^2}), which is currently an open problem [12, 21, 25].

Unfortunately, for demand matrices with high skew like those found in practice, BvN decomposition will often produce schedules with a large number of configurations. Many of these configurations will have short durations (e.g., on the order of the reconfiguration delay, δ), leading to low overall efficiency. The crux of the issue is that BvN must serve every demand eventually, which means that configurations with low duration (and thus low efficiency) must be used eventually.

Algorithm 2: Solstice

```

input : The demand:  $D$ , reconfiguration delay:  $\delta$ ,
         link rate for circuit switch:  $r_c$ ,
         link rate for packet switch:  $r_p$ 
output:  $m$  circuit configurations and durations:  $\{P_i\}, \{t_i\}$ ,
         demand sent to packet switch:  $E$ 
 $E \leftarrow D$ 
 $D' \leftarrow \text{QuickStuff}(D)$ 
 $T \leftarrow 0$ 
 $r \leftarrow$  largest power of 2 smaller than  $\max(D')$ 
 $i \leftarrow 1$ 
while  $\exists$  row or column sum of  $D' > r_p T$  do
   $P_i \leftarrow \text{BigSlice}(D', r)$ 
  if  $P_i \neq \text{NULL}$  then
     $t_i \leftarrow \min\{D'_{a,b} \mid P_{i,a,b} = 1\}/r_c$ 
     $D' \leftarrow D' - r_c t_i P_i$ 
     $E \leftarrow E - r_c t_i P_i$ 
     $E \leftarrow \text{ZeroEntriesBelow}(E, 0)$ 
     $T \leftarrow T + t_i + \delta$ 
     $i \leftarrow i + 1$ 
  else
     $r \leftarrow r/2$ 
  end
end
 $m \leftarrow i - 1$ 

```

If some demand can be ignored—or served by the packet switch in a hybrid network—then the focus of the scheduling algorithm shifts towards finding configurations that can be used for long time durations, thus, increasing efficiency.

4.2 Solstice algorithm

Our initial assumptions about the demand matrix (namely sparsity and skewness; see §2.3) motivate our scheduler’s design. Despite BvN’s potential to generate $O(n^2)$ configurations, a sparse demand matrix lowers this to $O(n)$ (assuming sparsity is bounded by a constant). This motivates using BvN as a basis for Solstice, as fewer configurations lead to less reconfiguration delay. The high skew of our demand matrices implies easier separability into “big” (circuit) and “small” (packet) demands in a hybrid environment, motivating a greedy heuristic to select configurations with large durations for the circuit switch.

We now present the Solstice hybrid scheduling algorithm, shown in Algorithm 2. Solstice consists of two stages: a round of *stuffing* followed by multiple iterations of *slicing*. Stuffing takes an arbitrary demand matrix D and adds artificial demand to compute D' which is k -bistochastic. Slicing builds on BvN, leveraging the decomposability of k -bistochastic matrices to iteratively compute a schedule of configurations with long durations, greedily avoiding short, inefficient configurations. Solstice terminates when the demand not yet satisfied in the current (iteratively computed) schedule is small enough to be satisfied by the packet switch.

In addition to computing the m circuit configurations $\{P_i\}$ and durations $\{t_i\}$, each iteration of Solstice’s slicing maintains three additional variables: r , T , and E , consistent with our notation in §2.2.1. r is a threshold for the current iteration (see §4.2.2), and T is the total time used by circuit configu-

Function QuickStuff

```

input : Demand matrix  $D$ 
output :  $k$ -bistochastic matrix  $D'$ 
 $D' \leftarrow D$ 
 $\{R_i\} \leftarrow$  The sums of each row in  $D'$ 
 $\{C_i\} \leftarrow$  The sums of each column in  $D'$ 
 $\phi \leftarrow \max(\{R_i\} \cup \{C_i\})$ .
for each  $D'_{i,j} > 0$  do
  | add  $\phi - \max(R_i, C_j)$  to  $D'_{i,j}$ ,  $R_i$ , and  $C_j$ 
end
for each  $D'_{i,j} = 0$  do
  | add  $\phi - \max(R_i, C_j)$  to  $D'_{i,j}$ ,  $R_i$ , and  $C_j$ 
end
Return  $D'$ .

```

rations computed so far. E is an $n \times n$ matrix encoding the amount of residual demand not serviced by the configurations computed so far. At termination, any demand left in E will be scheduled on the packet switch.

4.2.1 Stuffing

Stuffing is a heuristic pre-processing step that converts the demand matrix D into a k -bistochastic matrix D' by adding artificial demand. As explained in §4.1, the BvN theorem proves the existence of a simple decomposition (i.e., schedule) of a matrix as long as the matrix is k -bistochastic. A naive stuffing approach is to go through each $D_{i,j}$ in order and increase its value (“stuff” it) until either the sum of row i or the sum of column j reaches k . Iterating over the elements of D , eventually all row/column sums will be k , as required. This approach is suboptimal because the entries of D that are increased may be entirely entries that were zero originally, making the matrix less sparse (leading to more costly reconfigurations).

A better approach would be to stuff the largest elements of D first, as the artificial demand needed to stuff these elements would be proportionally smaller; this approach, however, is computationally expensive. Solstice, instead, uses the stuffing function is listed in Function QuickStuff. Instead of sorting all the elements, QuickStuff simply stuffs the non-zero elements of D in arbitrary order, providing a reasonable approximation. Afterwards, the zero elements are visited in order to stuff any elements that still need to be stuffed. Focusing on non-zero elements first helps (in the average case) keep the resulting matrix sparse. In practice, QuickStuff keeps the sparsity of D' similar to D , but we leave the theoretical analysis of its worst case to future work.

4.2.2 Slicing

After stuffing, Solstice enters its second stage, which is conceptually similar to BvN’s main loop: finding the next circuit configuration and its corresponding duration. We call this process slicing. Examining all possible configurations has no known polynomial-time algorithm; Solstice picks one greedily. Solstice differs from BvN in that it selects configurations with long corresponding durations to amortize the reconfiguration penalty and keep utilization high. Also,

Function BigSlice

```

input :  $k$ -bistochastic matrix  $D'$ , threshold  $r$ 
output : A permutation matrix  $P$  s.t.
   $\{D'_{a,b} \mid P_{a,b} = 1\} \geq r$ 
 $D'' \leftarrow$  ZeroEntriesBelow( $D'$ ,  $r$ )
 $B \leftarrow$  BinaryMatrixOf( $D''$ )
Interpret  $B$  as a bipartite graph of senders to receivers.
if  $\exists P$ , a perfect matching of  $B$  then
  | return  $P$  interpreted as a permutation matrix.
else return NULL ;

```

unlike BvN, Solstice terminates once unscheduled traffic can be feasibly forwarded by the packet switch.

Each iteration of slicing, shown in Function BigSlice, takes as input a (stuffed) matrix D' and a threshold r and returns a circuit configuration P_i such that when D' is overlayed onto the links in P_i , each link has at least r bits ready to send. If no such configuration P_i exists, NULL is returned. If we interpret D' as a bipartite graph of senders and receivers, Solstice effectively tries to find a Max-Min Weighted Matching (MMWM), which is the perfect matching (i.e., a matching of size n) with the largest minimum element. Like BvN, a particular circuit configuration’s duration is set based on the minimum element for that configuration.

The MMWM search is iterative: Slicing starts with a high threshold (r) of the largest power of 2 smaller than the maximum element in the stuffed demand matrix and then tries to find an arbitrary perfect matching on the stuffed demand matrix where values less than the threshold are ignored. Thus, any perfect matching returned will have corresponding duration at least r/r_c .

Solstice keeps looking for perfect matchings using the same threshold until there are no longer any perfect matchings with this threshold. At that point, the threshold is reduced by half. An optimal MMWM algorithm would consider all $O(\max(D'))$ different thresholds; Solstice considers an exponentially spaced set of them.

Slicing ends when the packet switch has enough capacity to handle the remaining demand. Solstice tracks the (so-far) unsatisfied demand using matrix E . Another variable T keeps track of the total time spent sending data over the circuit switch, as well as the time spent reconfiguring (i.e., $T = (\sum_{i=1}^m t_i) + m\delta$). Once the total time T is large enough that the packet switch can handle the leftover demand E , Solstice terminates. When the algorithm terminates, no link on the packet switch—row or column sums in E —is required to send more than $r_p T$. As stuffing never increases the max row/column sum, we know that the max row/column sum of E and D' are always the same, allowing us to use them interchangeably in the loop termination condition. Phrased differently, Solstice terminates when roughly (r_p/r_c) of the traffic is allocated to the packet switch.

4.3 Example

For clarity, we now describe how Solstice operates on the simple demand matrix in Figure 3, with $\delta = 1$, $r_p = 0.1$, $r_c = 1$. We define the diameter of a matrix, D_{diameter} as the

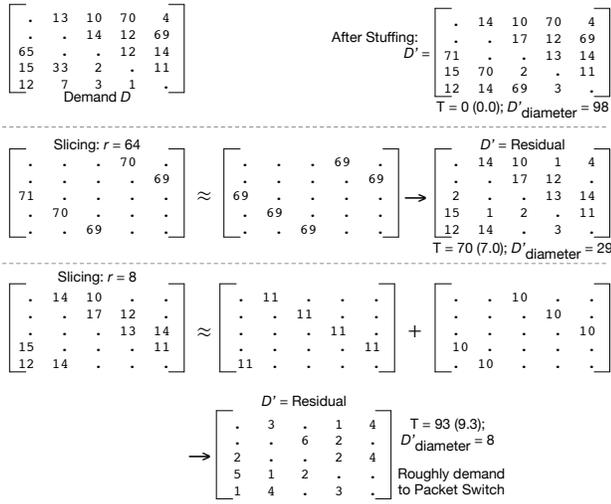


Figure 3: A sample execution with $\delta = 1$, $r_p = 0.1$, $r_c = 1$. Solstice computes a circuit schedule with 3 configurations with a total time of $T = 69 + 11 + 10 + 3\delta = 93$, and at most 8 bits of demand sent over each packet switch link. The values in parentheses show $r_p T$.

maximum row or column sum. The input demand matrix D has a diameter of 98 (the second row sum), so D is stuffed to obtain a matrix D' where each row and column sum is 98 (i.e., D' is k -bistochastic, with $k = 98$). In the first iteration, Solstice considers a minimum duration of $r = 64$ by extracting the subset of elements of at least that size. That subset admits only one perfect matching with a minimum-sized element of 69, so the duration of the first configuration is determined to be 69. It subtracts the demand from the stuffed matrix, D' . The total time is now $T = 69/r_c + 1\delta = 70$. If $r_p T$ (here 7) is greater than the diameter of D' (29), we could leave the rest to the packet switch. Unfortunately, this is not yet the case, so the algorithm continues to consider thresholds of exponentially decreasing size.

At $r = 32$ there are no elements considered, as none are larger than 32. At $r = 16$, one element is considered. As Solstice looks for perfect matching, at least n elements need to be considered, so the threshold is reduced again. Once $r = 8$, D' finally admits another perfect matching. This time, there are two: the first perfect matching BigSlice identifies has a minimum element of size 11, and the second matching it identifies has size 10. After accounting for the demand serviced by both of these configurations, the total time becomes $T = 69/r_c + 11/r_c + 10/r_c + 3\delta = 93$, at which point we can schedule the rest on the packet switch (as $8 \leq 9.3$), allowing Solstice to avoid scheduling the “long tail” of demands on the circuit switch, leading to high utilization. The actual demand sent to the packet switch, E , is slightly less than the residual D' , as D' is the result of stuffing. The actual values for E are omitted for brevity.

4.4 Worst case

For very sparse matrices, Solstice performs nearly optimally (see §5.5); however, a very sparse demand matrix fun-

damentally has fewer viable circuit configurations to choose from, simplifying scheduling. A very dense demand matrix may appear difficult to schedule, but is simple as well; an all-to-all workload can be scheduled efficiently using weighted round robin (i.e., n schedules). Solstice correctly identifies this, as we show in §5.5, and uses exactly n schedules.

The extremes of skew also may appear problematic, but are similarly straightforward to schedule. Demand matrices with very high skew reduce to the simple problem of identifying large flows (for the circuit switch) and small flows (for the packet switch). Very low skew is more efficiently solved through weighted round robin, similar to above.

The analysis of Solstice’s worst case (both in terms of schedule duration and number of configurations) is currently unknown, because of the difficulty in characterizing the types of workloads for which Solstice does poorly. The regions between these extremes, both for sparsity and skew, appear to not have simple solutions. Thus, we base our evaluation (§5) on these difficult regions. It is currently unclear how to mathematically characterize these regions exactly, so we leave the analysis of Solstice’s worst case as future work.

Finally, as Solstice always schedules all traffic in the demand matrix completely, each demand matrix over time is completely independent. Thus, inter-demand matrix patterns (e.g., alternating structure across matrices) are not a problem for Solstice.

4.5 Time complexity

QuickStuff runs in $O(n^2)$ time. Unfortunately, it is possible that QuickStuff will output a matrix with less sparsity, i.e., D'_{count} (the maximum number of non-zero elements in a row or column of D') may be greater than D_{count} , as stuffing may add to entries that were zero. In the general case, D'_{count} may approach n . For input matrices with low sparsity, however, our experience shows that it is rarely substantially larger, as QuickStuff focuses on stuffing non-zero entries.

BigSlice itself is dominated by the matching algorithm step. Goel et al. propose a randomized matching algorithm [9] that, in expectation, takes $O(|V| \log^2 |V|)$ plus a one-time preprocessing step of $O(|E|)$, or in our terms $O(n \log^2 n)$ and $O(nD'_{\text{count}})$. In the worst case $D'_{\text{count}} = n$, resulting in an initial $O(n^2)$ preprocessing step.

In the worst case, BigSlice will need to try $O(nD'_{\text{count}})$ different thresholds (i.e., $\max(D)$ is certainly less than $2^{nD'_{\text{count}}}$) and may additionally need $O(nD'_{\text{count}})$ successful calls to BigSlice that generate configurations (i.e., each schedule zeros only one element of the matrix), thus requiring $O(nD'_{\text{count}} + nD'_{\text{count}}) = O(nD'_{\text{count}})$ calls to BigSlice. In total, slicing takes $O(nD'_{\text{count}} + nD'_{\text{count}} * n \log^2 n) = O(D'_{\text{count}} n^2 \log^2 n)$, which falls into line with Goel’s analysis for BvN.

Stuffing takes $O(n^2)$ and slicing takes $O(D'_{\text{count}} n^2 \log^2 n)$; other smaller steps either take $O(1)$ or $O(nD'_{\text{count}})$ time. Thus, the overall complexity of Solstice is $O(D'_{\text{count}} n^2 \log^2 n)$. It is possible that D'_{count} can be n due to stuffing, leading to $O(n^3 \log^2 n)$. Empirically, we find D'_{count} is effectively constant, yielding $O(n^2 \log^2 n)$.

5. EVALUATION

We evaluate the performance of Solstice along three distinct dimensions to answer the following questions:

1. **Utilization:** How does Solstice perform compared to classic algorithms for switches of varying port count? (Solstice performs up to $2.9\times$ better than BvN.)
2. **Skew:** How does Solstice handle varying the weight of small demands? Is high skew required for Solstice to perform well? (Solstice performs only 10% worse with low skew compared to our baseline.)
3. **Sparsity:** How does Solstice handle very sparse and very dense matrices? Are sparse matrices required for Solstice to perform well? (Solstice performs only 26% worse in a completely filled matrix compared to our baseline.)

Additionally, we reflect on Solstice’s performance by considering whether there are simple ways to improve upon it and how it would function as a non-hybrid circuit scheduler.

5.1 Bounds

As the ILP presented in §3 is too slow to compute for port counts larger than ~ 12 , it is difficult to compare Solstice against the truly optimal schedule duration (T) for a given demand matrix D . However, it is possible to provide weak lower and upper bounds on the optimal T .

Lower bound: We mathematically derive a weak lower bound. As it is a lower bound, it is impossible to build a schedule that completes in less time. However, it is weak in the sense that it might not be possible to build a schedule that completes in that amount of time. We provide intuition for the lower bound by starting with a purely circuit switched network, LB_c :

$$LB_c = D_{\text{diameter}}/r_c + D_{\text{count}} * \delta.$$

Serving demand D on a pure circuit switch requires at least as much time as the largest row or column sum, D_{diameter} , divided by the link rate. It also requires at least as many configurations as the largest row or column count, D_{count} , each incurring a reconfiguration penalty of δ .

For a hybrid switch, we relax our assumption that we need D_{count} configurations (as in some cases the count of every row and column may be reduced to one or zero by sending data over the packet switch). Thus, we weaken the bound to only one reconfiguration penalty (δ). Moreover, with the addition of a packet switch, the total amount of time needed can be reduced proportionally:

$$LB_h = (r_c/(r_c + r_p)) * (D_{\text{diameter}}/r_c + \delta).$$

For example, let the link rate of the circuit switch $r_c = 10$, and the link rate of the packet switch $r_p = 1$. For every 11 units of demand that come in, 10 can be sent to the circuit switch and 1 can be sent to the packet switch, while overall taking 10 units of time. Thus, multiplying our previous lower bound by $r_c/(r_c + r_p) = 10/11$ accounts for the inclusion of the packet switch.

Upper bound: All correct scheduling algorithms provide upper bounds on the optimal T . Once we have computed a

Algorithm	Runtime
Lower bound	< 1 ms (64 ports)
BvN	27 ms (64 ports)
iSLIP	13 ms (64 ports)
Solstice	2.9 ms (64 ports)
Solstice++	5 min (64 ports)
Optimal	5 min (12 ports); $\gg 11$ hours (16 ports)

Table 3: The runtime of each scheduling approach on an Intel® Xeon® E5-2680 v2 (2.8 GHz) processor with 128 GB of memory.

schedule, we never need to consider any schedules with a longer duration (i.e., the schedule is an upper bound on all schedules that serve D). We optimize the schedule produced by Solstice (in a tractable but non-realtime manner) to provide a weak upper bound on the optimal T . It is weak in the sense that it provides a feasible solution, but not necessarily the best possible solution. It is, however, the best upper bound of which we are aware, as it always produces a schedule no worse than Solstice, and we are not aware of any algorithm that provides better schedules than Solstice for the hybrid scheduling problem.

We observe that schedules computed by Solstice can be optimized by using the resulting configurations (but ignoring their durations) as the candidate permutations for the ILP from §3. In addition, we enhance the candidate set with a small number of randomly generated permutations. The ILP will clearly produce a schedule no worse than that selected by Solstice, but frequently is able to compute better durations, “throw away” some steps in the schedule, and occasionally determine that one of the random permutations is a better choice (i.e., Solstice explored a local minimum in its search). Because this process involves solving an ILP, it is impractical for use as a scheduler. However, by iterating through this process several times (10 in our simulations), we can often improve Solstice. We call this scheduling algorithm Solstice++.

5.2 Simulation setup

Unless otherwise specified, our simulations consider a hybrid switch with 64 ports consisting of a 100-Gbps (per link) circuit switch with a reconfiguration time of $20 \mu\text{s}$ and a 10-Gbps (per link) packet switch (a 10:1 ratio). We consider scheduling 3 ms of demand at a time (as in ReactoR [19]) and assume demand matrices are sparse (4 large demands and 12 small demands per port) and skewed (small demands only make up 30% of total demand). We test sensitivity to each of these parameters in our evaluation.

We compare six different scheduling approaches: three practical scheduling algorithms—BvN (§ 4.1), Solstice, and an improved¹ version of iSLIP (a classical crossbar scheduling algorithm designed to be starvation free and easy to implement [20])—and three bounds: the optimal schedule computed (when tractable) by an ILP (§3), our lower bound on the optimal schedule duration, and an upper bound on the

¹iSLIP can produce schedules with repeated configurations. We merge all duplicate configurations into one with a longer duration.

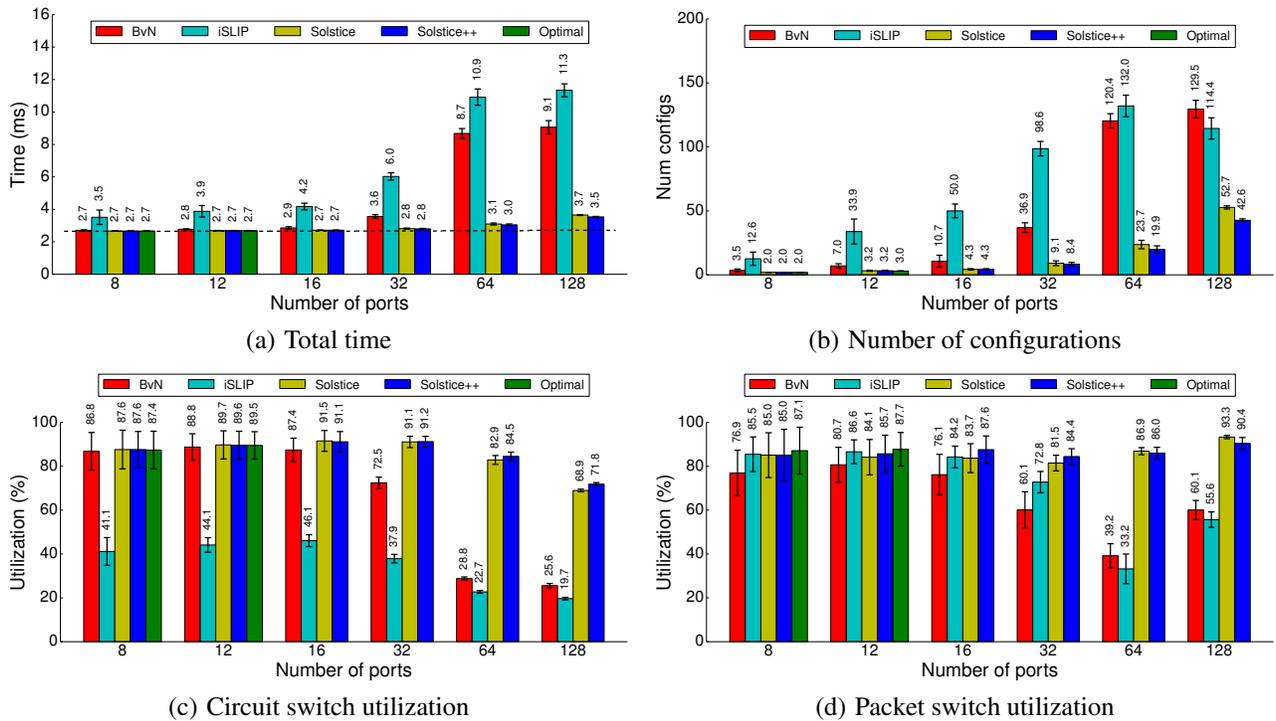


Figure 4: The performance of different scheduling approaches. Each bar represents the average of 100 runs; error bars show standard deviation. The dashed line represents a weak lower bound on optimal schedule duration. True optimal is somewhere between the dashed line and Solstice++.

optimal schedule duration computed by Solstice++. As both BvN and iSLIP are iterative algorithms, we iterate until the residual demand is small enough to be served by the packet switch. For context, example runtimes of each approach are listed in Table 3.

Demand: We construct traffic demand matrices based upon skew and sparsity characteristics derived from published data-center workloads, but with significantly higher overall traffic demands to stress the scheduler. We base the matrix generator on traces from the University of Wisconsin [2] and Alizadeh et al. [1]. The Wisconsin study provides two one-hour traces of traffic among 500–1000 servers. Examining the traffic matrices for each 3-ms window of traffic in this trace, the maximum number of non-zero elements in a matrix are just 36 and 85 (out of 500 and 1000 servers, respectively). The links for most hosts are mostly idle, and no host exchanges a large flow with more than five other hosts in a window. The Alizadeh work describes flow behavior and size distributions of a workload combining query traffic (one host sends to all other hosts) and background flows from other applications. Even when scaling this workload to be five times more dense, there are a maximum of seven concurrent large background flows per host in a 3-ms window (the paper shows at most four concurrent large flows per host in a 50-ms window). The small flow query traffic consumes about 10% of the switch capacity, and the large background flows use about 30%.

Constructing matrices from the demand model: To generate demand matrices that match the distributions above,

we generate workloads that have a fixed number of flows per source port. The default value is for 4 large flows and 12 small flows, which we vary in our evaluation of sparsity (§5.5). By default, for a given link, the large flows are given 70% of the link bandwidth (to split evenly) and the small flows are given 30% (to split evenly), which we vary in our evaluation of skew (§5.4). To avoid completely saturating each link, we scale the result back to 96% of the total link bandwidth. Finally, the demands are perturbed with noise by adding $\pm 0.3\%$ of the link bandwidth.

The destination of each flow is selected in one of two ways: controlled or random. By default, we construct demand matrices in a controlled fashion by assigning flow destinations through combining multiple randomly generated circuit configurations (i.e., permutation matrices). The resulting matrices are “controlled” as they have structure in their communication and closely match the workloads upon which we base our traffic models, but are also somewhat easier to schedule (the solution involves decomposing the demand back into the original circuit configuration, though this is not strictly possible due to the addition of noise). For comparison, we also evaluate workloads where destinations are assigned in a purely random fashion (§5.4.2).

5.3 Utilization

We begin by considering how effectively Solstice schedules demand—i.e., the utilization it is able to achieve—as switch port count increases from 8 to 128. We show the results (Figure 4) in terms of total time to satisfy all demand,

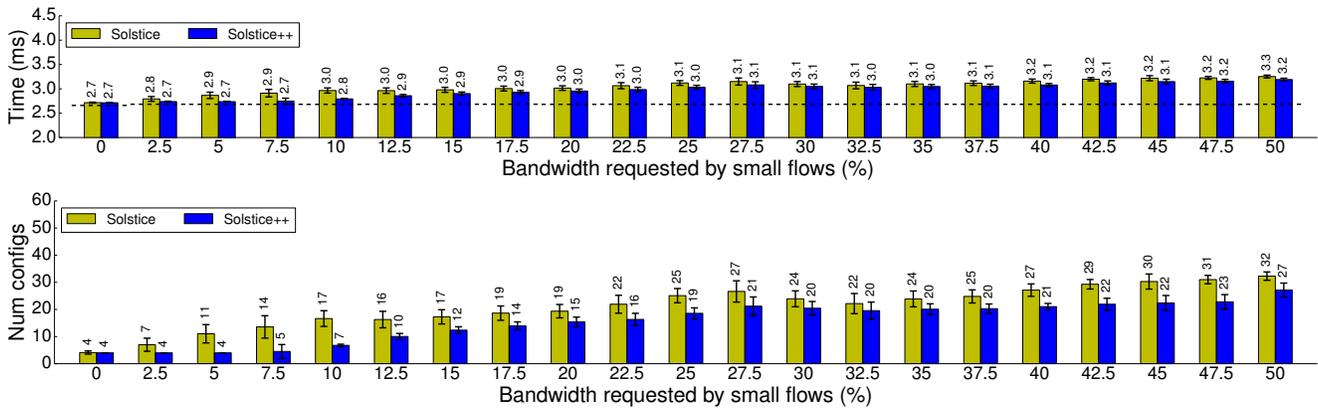


Figure 5: Total time (top) and number of configurations (bottom) as a function of controlled skew. Each bar represents the average of 25 runs; the error bars show standard deviations. Note the top y-axis starts at 2 ms. The dashed line represents a weak lower bound on optimal schedule duration. True optimal is somewhere between the dashed line and Solstice++.

the inverse of utilization. In order to understand why each approach performs as it does, we also plot the number of circuit configurations each includes in its schedule and the resulting utilization on both the circuit and packet switches.

Total time: Figure 4(a) plots the lower bound on optimal (LB_h) as a horizontal dashed line. We continue this convention throughout our evaluation. Recall that Solstice++ is an upper bound on optimal. This means that true optimal is somewhere between the dashed line and Solstice++.

For small scales (8 and 12 ports), we see that both Solstice and BvN achieve the lower bound, which turns out to be tight (Optimal performs no worse). We cannot compute Optimal at medium scales (16 and 32 ports), but Solstice is only slightly slower than the lower bound. Moreover, it performs similarly to Solstice++, suggesting that perhaps the lower bound is no longer tight. At larger scales (64 and 128 ports) we start to see divergence between Solstice and Solstice++. Despite further diverging from the lower bound, it is worth reiterating that the lower bound is loose; it is likely that Optimal would also diverge from the lower bound if it were tractable to compute at larger scales. For comparison, off-the-shelf algorithms BvN and iSLIP perform well until large scales, where they perform almost $3\times$ worse than Solstice.

Number of configurations: While Solstice performs as well as Optimal at small scale, the number of configurations starts to diverge at as early as 12 ports (Figure 4(b)). Similarly, Solstice uses more configurations than Solstice++ at 64 ports—but only slightly increases the total time due to the relatively small configuration penalty. The takeaway is Solstice includes a few redundant configurations (a point we will explore in §5.6), but the redundancy does not substantially impact the efficiency of the schedule.

Looking at how many configurations are produced by each algorithm provides insight into BvN and iSLIP’s poor performance. Both use a large number of configurations as they do not consider the reconfiguration penalty. We know of no mathematical assurance on how many configurations BvN produces in the average case. While iSLIP’s running time

(Table 3) is less than half of BvN, its often produces many more configurations.

Utilization Breakdown: We plot the utilization of both the circuit switch and packet switch in Figures 4(c) and 4(d), respectively. We see that Solstice performs similarly ($\sim 80\text{--}90\%$ circuit utilization) to Solstice++ (and Optimal, where available) between 8–64 ports. At our largest scale, 128 ports, circuit utilization drop precipitously. This is the direct result of the large increase in reconfigurations at this scale (Figure 4(b)): The 20% increase in the number of configurations generated by Solstice when compared to Solstice++ is reflected here as a 3% decrease in circuit utilization.

5.4 Skew

Solstice is designed to take advantage of skew in demand across flows, tailoring its heuristics to schedule workloads consisting of a small number of flows with (relatively) large demands among a background of many more flows with small demands. Here, we explore the behavior of Solstice under varying degrees of demand skew among a fixed number of flows. We explore skew for workloads where the destinations are selected by generating random circuit configurations (controlled) and where they are selected randomly.

5.4.1 Controlled skew

Figure 5 shows the total time (top) and number of circuit configurations (bottom) used as the skew changes. As expected, Solstice produces longer schedules when more traffic is used by the small flows, because the switch must be reconfigured more often to support them. Conversely, with little traffic in small flows, they can be completely tossed to the packet switch and the circuit switch only reconfigured for the number of large flows (i.e., 4).

Because the number and size of the demands stay the same across the experiments (modulo matrices where small and large demands overlap), the maximum row or column sum (D_{diameter}) and the maximum number of non-zero elements (D_{count}) stay the same, explaining the constancy in the lower bound (dashed line). Notably, Solstice++ deviates from the

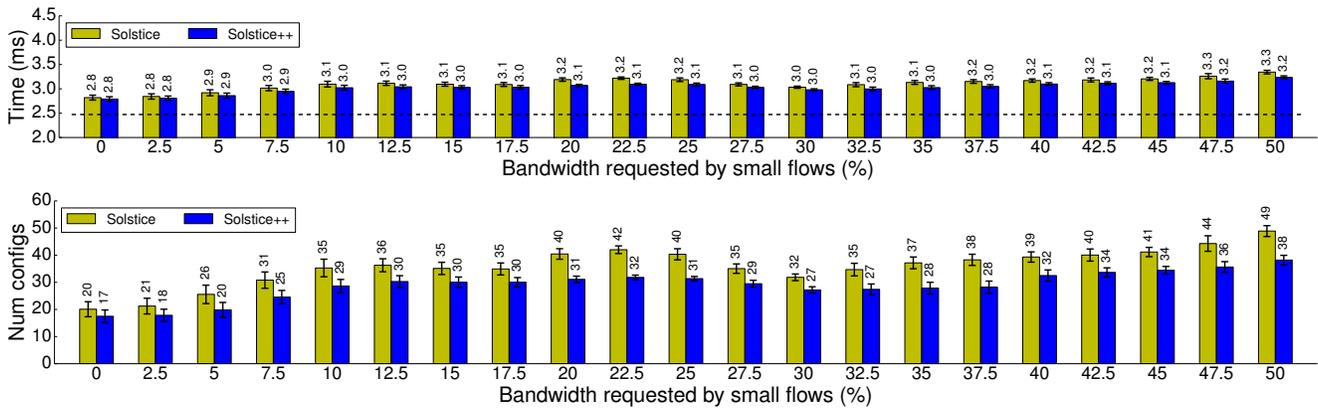


Figure 6: Total time (top) and number of configurations (bottom) as a function of random skew. Each bar represents the average of 25 runs; the error bars show standard deviations. Note the top y-axis starts at 2 ms. The dashed line represents a weak lower bound on optimal schedule duration. True optimal is somewhere between the dashed line and Solstice++.

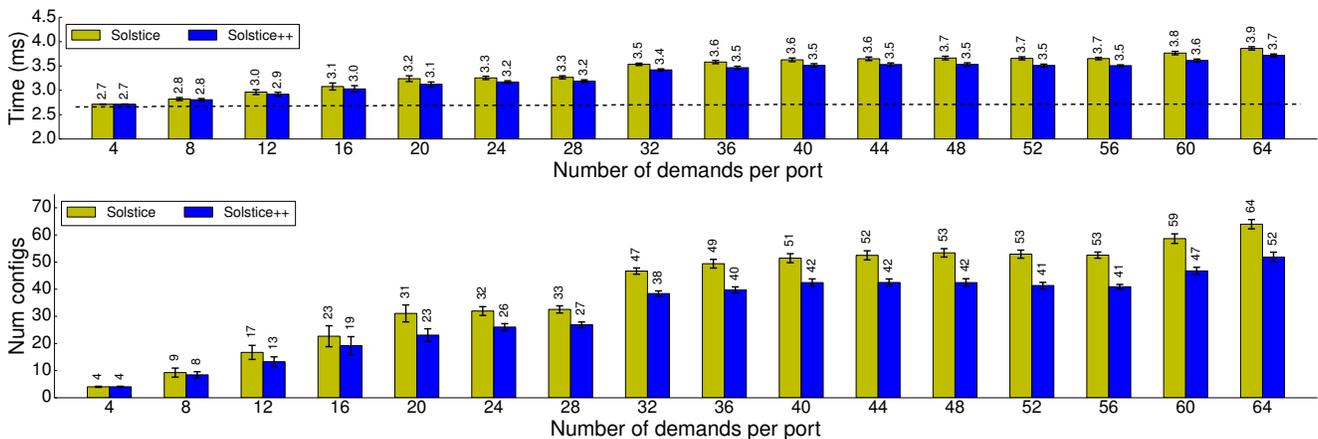


Figure 7: Total time (top) and number of configurations (bottom) as a function of sparsity. Each bar represents the average of 25 runs; the error bars show standard deviations. Note the top y-axis starts at 2 ms. The dashed line represents a weak lower bound on optimal schedule duration. True optimal is somewhere between the dashed line and Solstice++.

lower bound, suggesting that our lower bound may become more and more loose, which is intuitively correct; as the demand becomes less skewed, it is harder to remove traffic by tossing it to the packet switch. Effectively, less skewed demand is fundamentally less efficient to schedule.

5.4.2 Random skew

Unlike the previous experiment, here we randomize the source-destination pairs used to create the demand matrices rather than generating them from circuit configurations. Figure 6 shows the results in terms of time (top) and number of configurations (bottom). We see a similar gap between Solstice and Solstice++ in terms of the number of configurations, with commensurate increases in total time. We note again that the lower bound (dashed line) is constant as the diameter and count of the matrices remain constant.

Randomness affects the absolute magnitude of the number of configurations quite strikingly when compared to the controlled skew experiments—as much as $5\times$ the number of large demands per port when skew is high (i.e., when small

flows request no bandwidth). Random demand is harder to separate into a small number of circuit configurations, as it was not originally drawn from a set of circuit configurations. As more configurations are needed to solve these random demand matrices, more total time is needed.

The key takeaway from both skew experiments is that Solstice performs better (both absolutely and compared to our lower and upper bounds) when the demand matrix is skewed, but its performance does not decline drastically even with skew is minimized.

5.5 Sparsity

Here, we adjust the sparsity of the demand matrix (by varying the number of communicating pairs) to test Solstice’s sensitivity to sparsity. We construct a demand matrix that has k big flows and $l = 3k$ small flows per source and destination pair by randomly generating $k + l$ circuit configurations. We use the same calculation from §5.2 to allocate 30% of demand to small flows. We vary $k + l$ from 4 to 64 flows to reduce the sparsity of the demand matrix. At 64 flows the matrix

is completely filled with non-zero entries. Fundamentally, a dense matrix is easy to schedule for (i.e., use weighted round-robin), but may require many reconfigurations, increasing the total schedule time.

The results are shown in Figure 7, again as time (top) and number of configurations (bottom). As the matrix gets filled with more entries, more configurations are required. In the extreme case where the matrix is completely filled (64 demands), weighted round-robin (i.e., 64 configurations) is roughly the best solution, as shown in the graph. We again see that Solstice includes more configurations than necessary (as indicated by Solstice++) but their inclusion only marginally impacts the total time. We expect the number of configurations to be roughly in line with the average number of demands (modulo data sent over the packet switch), as we see in the graph. We see that the total time increases much faster than in the skew graphs, implying that Solstice is more sensitive to sparsity than skewness.

5.6 Discussion

In addition to the simulation results presented above, we consider additional extensions to Solstice but find they do not substantially improve its performance. We also explore whether Solstice is suitable as a traditional crossbar scheduler.

Improving Solstice: Solstice++ (as presented throughout the evaluation) improves Solstice’s results by considering multiple extra configurations and throwing away unnecessary ones. We find that the configurations Solstice++ deems necessary are almost always ($\geq 99.5\%$) a subset of the configurations Solstice used. Phrased differently, there is rarely benefit from considering configurations not identified by Solstice, which motivates exploring whether one could improve upon Solstice strictly by reconsidering how it employs the configurations it computes.

However, we find using an LP (similar to the ILP presented in §3) to adjust the time durations of Solstice’s configurations without removing any does not provide improvement. This implies that Solstice’s time selection is optimal. This makes sense as it always uses the minimum element of a slice as its time duration. QuickStuff minimally impacts scheduling, as durations picked by Solstice are based on the stuffed matrix, whereas the LP operates on the demand matrix directly.

Solstice on purely circuit networks: We re-run the utilization experiments from Section 5.3 without the packet switch—in other words, we consider Solstice’s performance as a traditional crossbar scheduling algorithm. The results are summarized in Table 4. Solstice performs much worse in such an environment as it can no longer move small “long tail” demands to the packet switch. Despite this, Solstice++ is able to perform much better by reducing the number of configurations by $\sim 77\%$, or $\sim 5\text{ms}$ of reconfiguration time. Solstice++ manages to do this by using longer, but more inefficient, durations for some schedules, as the benefit of avoiding additional configurations to clean up the tail greatly outweighs the inefficiency. The gap between Solstice++ and the lower bound also grows, hinting that the lower bound may be very loose.

Algorithm	12 port	64 port	128 port
Lower bound	2.96	3.25	3.60
Solstice	3.20 (15.40)	6.54 (180.51)	9.56 (329.93)
Solstice++	2.97 (3.21)	3.93 (34.53)	5.11 (75.80)
Optimal	2.96 (3.01)	-	-

Table 4: Performance for a purely circuit switch. Presented as total time (in ms) followed by number of configurations in parentheses. Each entry corresponds to an average over 100 runs.

6. RELATED WORK

The crossbar switch scheduling problem has been studied for decades. The basic approach, often referred to as time slot assignment (TSA), decomposes an accumulated demand matrix into a set of weighted permutation matrices. Classical results [3] and early work on scheduling satellite-switched time-division multiple access (SS/TDMA) systems [14] show how to compute a perfect schedule, but the resulting schedules consist of $O(n^2)$ configurations. Although this approach is optimal for a switch with trivial reconfiguration time, it performs poorly in our network model.

On the opposite end of the spectrum, when reconfiguration time is large, there exist algorithms [10, 26, 29] that use the fewest possible number of configurations (n). For moderate reconfiguration times, DOUBLE [26] computes a schedule that requires twice the minimum number of configurations, $2n$. Further improved algorithms [7, 18, 30] take the actual reconfiguration delay into account. These algorithms, however, do not benefit from sparse demand matrices, continuing to require $O(n)$ configurations to cover the demand.

Other existing work uses a speedup factor (i.e., the ratio of the internal transfer rate to the port link rate). Perhaps the most well known example is iSLIP [20], which requires a $2\times$ speedup to maintain stability. Many of these algorithms perform poorly (i.e., introduce large delays) when the traffic demand is skewed, leading others to suggest using randomization to address the issue [8].

7. CONCLUSION

The ever-increasing demand for low-cost, high-performance network fabrics in datacenter environments has generated tremendous interest in alternative switching architectures. Researchers have proposed hybrid switches that combine circuit and packet switching technologies but have stopped short of addressing scheduling. We take the first steps by characterizing the problem, exploring the space of possible scheduling algorithms, and glean insights based on their results. We craft an algorithm, Solstice, that takes advantage of sparsity and skewness observed in real datacenter traffic to provide $2.9\times$ higher circuit utilization when compared to traditional schedulers in hybrid environments, while being within 14% of optimal, at scale.

Our evaluation of scheduling algorithms sheds light on the challenges of scheduling for both hybrid and pure circuit networks. The performance gained by both Solstice and the ILP-assisted formulations over traditional schedulers is the result of the insight that inefficient short duration “tail”

configurations of traditional pure circuit schedules can be efficiently handled by a packet switch. We believe this insight can lead to the development of heuristic approximation algorithms for the pure circuit case, which might leverage indirection or careful cluster scheduling to avoid the need for expensive n -to- n connectivity for small flows. These ideas bear further theoretical and practical examination.

Acknowledgments

The authors would like to thank the National Science Foundation (NSF CNS-1314921 and CNS-1314721), Google (Google Focused Research Award), Microsoft Research, and Intel via the Intel Science and Technology Center for Cloud Computing (ISTC-CC) for their support. Additionally, the authors would like to thank Daniel M. Kane and Russell Impagliazzo for their insight regarding crossbar scheduling, and the anonymous reviewers for their feedback.

8. REFERENCES

- [1] M. Alizadeh, A. Greenberg, D. A. Maltz, J. Padhye, P. Patel, B. Prabhakar, S. Sengupta, and M. Sridharan. Data Center TCP (DCTCP). In *Proc. ACM SIGCOMM*, Aug. 2010.
- [2] T. Benson, A. Akella, and D. A. Maltz. Network Traffic Characteristics of Data Centers in the Wild. In *Proc. ACM IMC*, Nov. 2010.
- [3] G. Birkhoff. Tres Observaciones Sobre el Algebra Lineal. *Univ. Nac. Tucumán Rev. Ser. A*, 1946.
- [4] K. Chen, A. Singla, A. Singh, K. Ramachandran, L. Xu, Y. Zhang, and X. Wen. OSA: An Optical Switching Architecture for Data Center Networks and Unprecedented Flexibility. In *Proc. USENIX NSDI*, Apr. 2012.
- [5] N. Farrington, G. Porter, Y. Fainman, G. Papen, and A. Vahdat. Hunting Mice with Microsecond Circuit Switches. In *Proc. ACM HotNets-XI*, Oct. 2012.
- [6] N. Farrington, G. Porter, S. Radhakrishnan, H. Bazzaz, V. Subramanya, Y. Fainman, G. Papen, and A. Vahdat. Helios: A Hybrid Electrical/Optical Switch Architecture for Modular Data Centers. In *Proc. ACM SIGCOMM*, Aug. 2010.
- [7] S. Fu, B. Wu, X. Jiang, A. Pattavina, L. Zhang, and S. Xu. Cost and Delay Tradeoff in Three-Stage Switch Architecture for Data Center Networks. In *Proc. IEEE High Perf. Switching and Routing*, July 2013.
- [8] P. Giaccone, B. Prabhakar, and D. Shah. Randomized Scheduling Algorithms for High-Aggregate Bandwidth Switches. *IEEE J. Sel. Areas in Comms.*, May 2003.
- [9] A. Goel, M. Kapralov, and S. Khanna. Perfect Matchings in $O(n \log n)$ Time in Regular Bipartite Graphs. In *ACM STOC*, June 2013.
- [10] I. S. Gopal and C. K. Wong. Minimizing the Number of Switchings in a SS/TDMA System. *IEEE Trans. Comms.*, June 1985.
- [11] Gurobi. Gurobi Optimization. <http://www.gurobi.com/>.
- [12] J. Haglund and J. Remmel. Rook Theory for Perfect Matchings. *Advances in Applied Math.*, Aug. 2001.
- [13] D. Halperin, S. Kandula, J. Padhye, P. Bahl, and D. Wetherall. Augmenting Data Center Networks with Multi-gigabit Wireless Links. In *Proc. ACM SIGCOMM*, Aug. 2011.
- [14] T. Inukai. An Efficient SS/TDMA Time Slot Assignment Algorithm. *IEEE Trans. Comms.*, Oct. 1979.
- [15] S. Kandula, J. Padhye, and P. Bahl. Flyways To De-Congest Data Center Networks. In *Proc. ACM HotNets-VIII*, Oct. 2009.
- [16] S. Kandula, S. Sengupta, A. Greenberg, P. Patel, and R. Chaiken. The Nature of Data Center Traffic: Measurements & Analysis. In *Proc. ACM IMC*, Nov. 2009.
- [17] R. Kapoor, A. C. Snoeren, G. M. Voelker, and G. Porter. Bullet Trains: A Study of NIC Burst Behavior at Microsecond Timescales. In *Proc. ACM CoNEXT*, Dec. 2013.
- [18] X. Li and M. Hamdi. On Scheduling Optical Packet Switches with Reconfiguration Delay. *IEEE JSAC*, Sept. 2003.
- [19] H. Liu, F. Lu, A. Forencich, R. Kapoor, M. Tewari, G. M. Voelker, G. Papen, A. C. Snoeren, and G. Porter. Circuit Switching Under the Radar with REACToR. In *Proc. USENIX NSDI*, Apr. 2014.
- [20] N. McKeown. The iSLIP Scheduling Algorithm for Input-Queued Switches. *IEEE Trans. Networking*, Apr. 1999.
- [21] Q.-K. Pan and R. Ruiz. A Comprehensive Review and Evaluation of Permutation Flowshop Heuristics to Minimize Flowtime. *Comp. & Op. Research*, Jan. 2013.
- [22] G. Porter, R. Strong, N. Farrington, A. Forencich, P.-C. Sun, T. Rosing, Y. Fainman, G. Papen, and A. Vahdat. Integrating Microsecond Circuit Switching into the Data Center. In *Proc. ACM SIGCOMM*, Aug. 2013.
- [23] A. Roy, H. Zeng, J. Bagga, G. Porter, and A. C. Snoeren. Inside the Social Network's (Datacenter) Network. In *Proc. ACM SIGCOMM*, Aug. 2015.
- [24] R. Sinkhorn and K. Paul. Concerning Nonnegative Matrices and Doubly Stochastic Matrices. *Pacific J. Math.*, May 1967.
- [25] M. Tandon, P. Cummings, and M. LeVan. Flowshop Sequencing with Non-Permutation Schedules. *Comp. & Chem. Eng.*, Aug. 1991.
- [26] B. Towles and W. J. Dally. Guaranteed Scheduling for Switches with Configuration Overhead. *IEEE Trans. Networking*, Oct. 2003.
- [27] G. Wang, D. G. Andersen, M. Kaminsky, M. Kozuch, T. S. E. Ng, K. Papagiannaki, M. Glick, and L. Mummert. Your Data Center Is a Router: The Case for Reconfigurable Optical Circuit Switched Paths. In *Proc. ACM HotNets-VIII*, Oct. 2009.
- [28] G. Wang, D. G. Andersen, M. Kaminsky, K. Papagiannaki, T. S. E. Ng, M. Kozuch, and M. Ryan. c-Through: Part-time Optics in Data Centers. In *Proc. ACM SIGCOMM*, Aug. 2010.
- [29] B. Wu and K. L. Yeung. Minimum Delay Scheduling in Scalable Hybrid Electronic/Optical Packet Switches. In *IEEE GLOBECOM*, Nov. 2006.
- [30] B. Wu, K. L. Yeung, and X. Wang. Improving Scheduling Efficiency for High-Speed Routers with Optical Switch Fabrics. In *IEEE GLOBECOM*, Nov. 2006.
- [31] X. Zhou, Z. Zhang, Y. Zhu, Y. Li, S. Kumar, A. Vahdat, B. Y. Zhao, and H. Zheng. Mirror Mirror on the Ceiling: Flexible Wireless Links for Data Centers. In *Proc. ACM SIGCOMM*, Aug. 2012.