

Explicit Path Control in Commodity Data Centers: Design and Applications

Shuihai Hu¹ Kai Chen¹ Haitao Wu² Wei Bai¹ Chang Lan³
Hao Wang¹ Hongze Zhao⁴ Chuanxiong Guo²
¹HKUST ²Microsoft ³UC Berkeley ⁴Duke University

Abstract

Many data center network (DCN) applications require explicit routing path control over the underlying topologies. This paper introduces XPath, a simple, practical and readily-deployable way to implement explicit path control, using existing commodity switches. At its core, XPath explicitly identifies an end-to-end path with a path ID and leverages a two-step compression algorithm to pre-install all the desired paths into IP TCAM tables of commodity switches. Our evaluation and implementation show that XPath scales to large DCNs and is readily-deployable. Furthermore, on our testbed, we integrate XPath into 5 applications to showcase its utility.

1 Introduction

Driven by modern Internet applications and cloud computing, data centers are being built around the world. To obtain high bandwidth and achieve fault tolerance, data center networks (DCNs) are often designed with multiple paths between any two nodes [3, 4, 13, 16, 17, 30], and Equal Cost Multi-Path routing (ECMP) [22], is the state-of-the-art for multi-path routing and load-balancing in DCNs [5, 16, 30].

In ECMP, a switch locally decides the next hop from multiple equal cost paths by calculating a hash value, typically from the source and destination IP addresses and transport port numbers. Applications therefore cannot explicitly control the routing path in DCNs.

However, many emerging DCN applications such as provisioned IOPS, fine-grained flow scheduling, bandwidth guarantee, etc. [5, 7, 8, 18, 20, 21, 24, 25, 39, 45], all require explicit routing path control over the underlying topologies (§2).

Many approaches such as source routing [36], MPLS [35], ATM [9], and OpenFlow [28] can enforce explicit path control. However, source routing is not supported in the hardware of the data center switches, which typically only support destination IP based routing. MPLS and ATM need a signaling protocol (e.g., Label Distribution Protocol for MPLS) to establish the path or virtual circuit. Hence they are typically used only for traffic engineering in core networks instead of application-level or flow-level path control. OpenFlow

in theory can establish fine-grained routing path by installing flow entries in the switches via the OpenFlow controller. But in practice, there are practical challenges such as limited flow table size and dynamic flow setup that need to be addressed (§6).

In order to address the scalability and deployment challenges faced by the above mentioned approaches, this paper presents XPath for flow-level explicit path control. XPath addresses the dynamic path setup challenge by giving a positive answer to the following question: can we pre-install all desired routing paths between any two servers? Further, XPath shows that we can pre-install all these paths using the destination IP based forwarding TCAM tables of commodity switches¹.

We note that one cannot enumerate all possible paths in a DCN as the number could be extremely large. However, we observe that DCNs (e.g., [2–4, 16, 17, 19]) do not intend to use *all* possible paths but a set of *desired* paths that are sufficient to exploit topology redundancy (§2.2). Based on this observation, XPath focuses on pre-installing these desired paths in this paper. Even though, the challenge is that the sheer number of desired paths in large DCNs is still overwhelming, e.g., a Fattree ($k=64$) has over 2^{32} paths among ToRs, exceeding the size of IP table with 144K entries, by many magnitudes.

To tackle the above challenge, we introduce a two-step compression algorithm, i.e., paths to path sets aggregation and path ID assignment for prefix aggregation, which is capable of compressing a large number of paths to a practical number of routing entries for commodity switches (§3).

To show XPath’s scalability, we evaluate it on various well-known DCNs with even millions of nodes (§3.3). Our results suggest that XPath effectively expresses tens of billions of paths using only tens of thousands of routing entries. For example, for Fattree(64), we pre-install 4 billion paths using 64K entries²; for HyperX(4,16,100), we pre-install 17 billion paths using 36K entries. With such algorithm, XPath easily pre-installs all desired paths

¹The recent advances in switching chip technology make it ready to support 144K entries in IP LPM (Longest Prefix Match) tables of commodity switches (e.g., [1, 23]).

²The largest routing entries among all switches.

into IP tables with 144K entries, while still reserving space to accommodate more paths.

To demonstrate XPath’s deployability, we implement it on both Windows and Linux platforms under the umbrella of SDN, and deploy it on a 3-layer Fattree testbed with 54 servers (§4). Our experience shows that XPath can be readily implemented with existing commodity switches. Through basic experiments, we show that XPath handles failure smoothly.

To showcase XPath’s utility, we integrate it into 5 applications (from provisioned IOPS [24] to Map-reduce) to enable explicit path control and show that XPath directly benefits them (§5). For example, for provisioned IOPS, we use XPath to arrange explicit path with necessary bandwidth to guarantee the IOPS provisioned. For network update, we show that XPath easily assists networks to accomplish switch update at zero traffic loss. For Map-reduce data shuffle, we use XPath to identify non-contention parallel paths in accord with the many-to-many shuffle pattern, reducing the shuffle time by over $3\times$ compared to ECMP.

In a nutshell, XPath’s primary contribution is that it provides a practical, readily-deployable way to pre-install all the desired routing paths between any s-d pairs using existing commodity switches, so that applications only need to choose which path to use without worrying about how to setup the path, and/or the time cost or overhead of setting up the path.

To access XPath implementation code, please visit: <http://sing.cse.ust.hk/projects/XPath>.

The rest of the paper is organized as follows. §2 overviews XPath. §3 elaborates XPath algorithm and evaluates its scalability. §4 implements XPath and performs basic experiments. §5 integrates XPath into applications. §6 discusses the related work and §7 concludes the paper.

2 Motivation and Overview

2.1 The need for explicit path control

Case #1: Provisioned IOPS: IOPS are input/output operations per second. Provisioned IOPS are designed to deliver predictable, high performance for I/O intensive workloads, such as database applications, that rely on consistent and fast response times. Amazon EBS provisioned IOPS storage was recently launched to ensure that disk resources are available whenever you need them regardless of other customer activity [24, 33]. In order to ensure provisioned IOPS, there is a need for necessary bandwidth over the network. Explicit path control is required for choosing an explicit path that can provide such necessary bandwidth (§5.1).

Case #2: Flow scheduling: Data center networks are built with multiple paths [4, 16]. To use such multi-

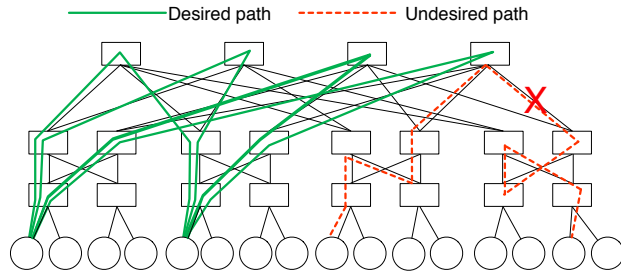


Figure 1: Example of the desired paths between two servers/ToRs in a 4-radix Fattree [4] topology.

ple paths, state-of-the-art forwarding in enterprise and data center environments uses ECMP to statically stripe flows across available paths using flow hashing. Because ECMP does not account for either current network utilization or flow size, it can waste over 50% of network bisection bandwidth [5]. Thus, to fully utilize network bisection, we need to schedule elephant flows across parallel paths to avoid contention as in [5]. Explicit path control is required to enable such fine-grained flow scheduling, which benefits data intensive applications such as Map-reduce (§5.5).

Case #3: Virtual network embedding: In cloud computing, virtual data center (VDC) with bandwidth guarantees is an appealing model for cloud tenants due to its performance predictability in shared environments [7, 18, 45]. To accurately enforce such VDC abstraction over the physical topology with constrained bandwidth, one should be able to explicitly dictate which path to use in order to efficiently allocate and manage the bandwidth on each path (§5.3).

Besides the above applications, the need for explicit path control have permeated almost every corner of data center designs and applications, from traffic engineering (*e.g.*, [8, 21]), energy-efficiency (*e.g.*, [20]), network virtualization (*e.g.*, [7, 18, 45]), to network congestion-free update (*e.g.*, [25]), and so on. In §5, we study 5 of them.

2.2 XPath overview

As introduced, XPath attempts to pre-install all desired paths using IP tables of commodity switches, so that applications can use these pre-installed explicit paths conveniently without any concern of dynamically setting up them. Next, we introduce what are the desired paths, and then overview the XPath framework.

Desired paths: XPath does not try to pre-install all possible paths in a DCN because this is impossible and impractical. We observe that when designing DCNs, people do not intend to use all possible paths in the routing. Instead, they use a set of *desired* paths which are sufficient to exploit the topology redundancy. This is the case for many recent DCN designs such as [2–4, 16, 17, 19]. For example, when designing a k -radix Fattree [4], they explicitly exploit $k^2/4$ parallel paths between any two ToRs

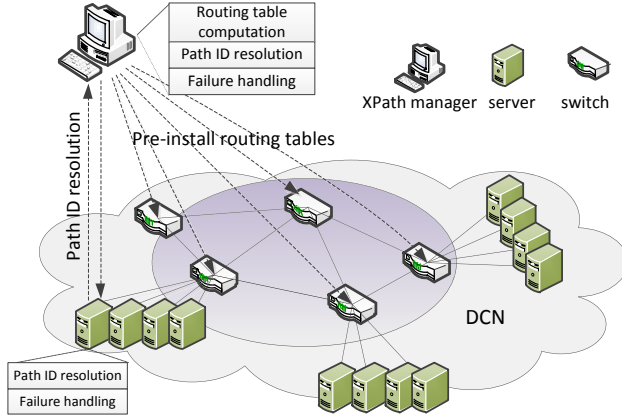


Figure 2: The XPath system framework.

for routing (see Fig. 1 for example); when designing an n -layer BCube [17], they use $(n + 1)$ parallel paths between any two servers. These sets of desired paths have already contained sufficient parallel paths between any s-d pairs to ensure good load-balancing and handle failures. Therefore, as the first step, XPath focuses on pre-installing all these desired paths.

XPath framework: Fig. 2 overviews XPath. As many prior DCN designs [11, 16, 17, 30, 40], XPath employs a logically centralized controller, called *XPath manager*, to handle all intelligence. The XPath manager has 3 main modules: routing table computation, path ID resolution, and failure handling. Servers have client modules for path ID resolution and failure handling.

- *Routing table computation:* This module is the heart of XPath. The problem is how to compress a large number of desired paths (e.g., tens of billions) into IP tables with 144K entries. To this end, we design a two-step compression algorithm: paths to path sets aggregation (in order to reduce unique path IDs) and ID assignment for prefix aggregation (in order to reduce IP prefix based routing entries). We elaborate the algorithm and evaluate its scalability in §3.
- *Path ID resolution:* In XPath, path IDs (in the format of 32-bit IP, or called routing IPs³) are used for routing to a destination, whereas the server has its own IP for applications. This entails path ID resolution which translates application IPs to path IDs. For this, the XPath manager maintains an IP-to-ID mapping table. Before a new communication, the source sends a request to the XPath manager resolving the path IDs to the destination based on its application IP. XPath may return multiple path IDs in response, providing multiple paths to the destination for the source to select. These path IDs will be locally cached for subsequent communications, but need to be forgotten periodically for failure handling. We elaborate this module and its implementation in §4.1.

³We use routing IPs and path IDs interchangeably in this paper.

- *Failure handling:* Upon a link failure, the detecting switch will inform the XPath manager. Then the XPath manager will identify the affected paths and update the IP-to-ID table (i.e., disable the affected paths) to ensure that it will not return a failed path to a source that performs path ID resolution. Similarly, upon a link recovery, the recovered paths will be added back to the IP-to-ID table accordingly. We rely on path ID swapping to handle failures. Because the source has cached multiple path IDs for a destination, once a link fails, we route around it by swapping the path IDs. In the meanwhile, the source will request the updated path IDs to the destination. After a failed link recovers, the source can use the recovered link once the local cache expires and a new path ID resolution is performed (§4.2).

Caveat: We note that XPath is expressive and is able to pre-install all desired paths in large DCNs into commodity switch IP tables, thus XPath’s routing table recomputation is performed infrequently, and cases such as link failures or switch upgrade [25] are handled through path ID swapping rather than switch table reconfiguration. However, table recomputation is necessary for extreme cases like network wide expansion where the network topology has fundamentally changed.

3 XPath Algorithm and Scalability

We elaborate the XPath two-step compression algorithm in §3.1 and §3.2. Then, we evaluate it on various large DCNs to show XPath’s scalability in §3.3.

3.1 Paths to path sets aggregation (Step I)

The number of desired paths is large. For example, Fat-tree(64) has over 2^{32} paths between ToRs, more than what a 32-bit IP/ID can express. To reduce the number of unique IDs, we aggregate the paths that can share the same ID without causing routing ambiguity into a non-conflict path set, identified by a unique ID.

Then, what kinds of paths can be aggregated? Without loss of generality, two paths have three basic relations between each other, i.e., convergent, disjoint, and divergent as shown in Fig. 3. Convergent and disjoint paths can be aggregated using the same ID, while divergent paths cannot. The reason is straightforward: suppose two paths diverge from each other at a specific switch and they have the same ID $path1 = path2 = path_id$, then there will be two entries in the routing table: $path_id \rightarrow port_x$ and $path_id \rightarrow port_y$, ($x \neq y$). This clearly leads to ambiguity. Two paths can be aggregated without conflict if they do not cause any routing ambiguity on any switch when sharing the same ID.

Problem 1: Given the desired paths $P = \{p_1, \dots, p_n\}$ of a DCN, aggregate the paths into non-conflict path sets so that the number of sets is minimized.

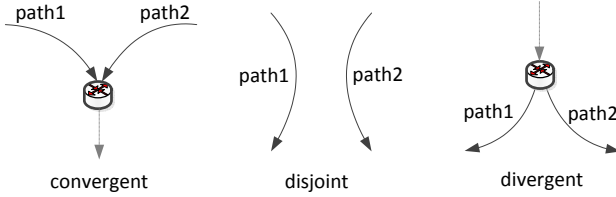


Figure 3: Three basic relations between two paths.

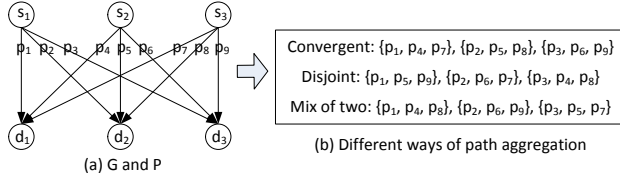


Figure 4: Example of different ways of path aggregation.

We find that the general problem of paths to non-conflict path sets aggregation is NP-hard since it can be reduced from the Graph vertex-coloring problem [41]. Thus, we resort to practical heuristics.

Based on the relations in Fig. 3, we can aggregate the convergent paths, the disjoint paths, or the mix into a non-conflict path set as shown in Fig. 4. Following this, we introduce two basic approaches for paths to path sets aggregation: convergent paths first approach (CPF) and disjoint paths first approach (DPF). The idea is simple. In CPF, we prefer to aggregate the convergent paths into the path set first until no more convergent path can be added in; Then we can add the disjoint paths, if exist, into the path set until no more paths can be added in. In DPF, we prefer to aggregate the disjoint paths into the path set first and add the convergent ones, if exist, at the end.

Actually, the obtained CPF path sets and DPF path sets have their own benefits. For example, a CPF path set facilitates many-to-one communication for data aggregation because such an ID naturally defines a many-to-one communication channel. A DPF path set, on the other hand, identifies parallel paths between two groups of nodes, such an ID naturally identifies a many-to-many communication channel for data shuffle.

In practice, users may have their own preferences to define customized path sets for different purposes. For example, BCube [17] identifies multiple edge-disjoint spanning trees in its topology for all-to-one traffic, and each of the spanning trees can be a customized path set. SecondNet [18] defines specific path sets for tenants according to tenant requests. Our XPath framework accommodates all forms of customized path sets as long as the path sets themselves are free of routing ambiguity.

3.2 ID assignment for prefix aggregation (Step II)

While unique IDs can be significantly reduced via Step I, the absolute value is still large. For example, Fattree(64)

Path set	Egress port	ID assignment (bad)	ID assignment (good)
$pathset_0$	0	0	4
$pathset_1$	1	1	0
$pathset_2$	2	2	2
$pathset_3$	0	3	5
$pathset_4$	1	4	1
$pathset_5$	2	5	3
$pathset_6$	0	6	6
$pathset_7$	0	7	7

Table 1: Path set ID assignment.

Path set	ID	Prefix	Egress port
$pathset_{1,4}$	0, 1	00*	1
$pathset_{2,5}$	2, 3	01*	2
$pathset_{0,3,6,7}$	4, 5, 6, 7	1**	0

Table 2: Compressed table via ID prefix aggregation.

has over 2 million IDs after Step I. We cannot allocate one entry per ID flatly with 144K entries. To address this problem, we further reduce routing entries using ID prefix aggregation. The insight is that a DCN is under centralized control and the IDs of paths can be coordinately assigned. Our goal is to assign (or order) the path IDs in such a way that they can be best aggregated using prefixes among all switches.

3.2.1 Problem description

We assign the IDs of paths that traverse the same egress port consecutively so that these IDs can be expressed using one entry through prefix aggregation. For example, in Table 1, 8 path sets go through a switch with 3 ports. A naive (bad) assignment will lead to uncompressed routing table with 7 entries. However, if we assign the paths that traverse the same egress port with consecutive IDs (good), we can obtain a compressed table with 3 entries as shown in Table 2.

To optimize for a single switch, we can easily achieve the optimal by grouping the path sets according to the egress ports and encoding them consecutively. In this way, the number of entries is equal to the number of ports. However, we optimize for all the switches simultaneously instead of one.

Problem 2. Let $T = \{t_1, t_2, \dots, t_{|T|}\}$ be the path sets after solving Problem 1. Assigning (or ordering) the IDs for these path sets so that, after performing ID prefix aggregation, the largest number of routing entries among all switches is minimized.

In a switch, a block of consecutive path IDs with the same egress interface can be aggregated into one entry⁴. We call this an aggregateable ID block (AIB). The number of such AIB s indicates routing states in the switch.

⁴The consecutiveness has local significance. Suppose path IDs 4, 6, 7 are on the switch (all exit through port p), but 5 are not present, then 4, 6, 7 are still consecutive and can be aggregated as $1** \rightarrow p$.

Thus, our goal is to minimize \mathcal{AIB} s among all switches through coordinated ID assignment.

To formally describe the problem, we use a matrix \mathbf{M} to describe the relation between path sets and switches. Suppose switches have k ports (numbered as $1\dots k$), then we use $m_{ij} \in [0, k]$ ($1 \leq i \leq |S|, 1 \leq j \leq |T|$) to indicate whether t_j passes through switch s_i , and if yes, what the egress port is. When $1 \leq m_{ij} \leq k$, it means t_j goes through s_i and the egress port is m_{ij} , and 0 otherwise means t_j does not appear on switch s_i .

$$\mathbf{M} = \begin{matrix} & t_1 & t_2 & t_3 & \dots & t_{|T|} \\ \begin{matrix} s_1 \\ s_2 \\ s_3 \\ \vdots \\ s_{|S|} \end{matrix} & \begin{pmatrix} m_{11} & m_{12} & m_{13} & \dots & m_{1|T|} \\ m_{21} & m_{22} & m_{23} & \dots & m_{2|T|} \\ m_{31} & m_{32} & m_{33} & \dots & m_{3|T|} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ m_{|S|1} & m_{|S|2} & m_{|S|3} & \dots & m_{|S||T|} \end{pmatrix} \end{matrix}$$

Then, we assign the IDs for path sets t_j ($1 \leq j \leq |T|$). We use $f(t_j) = r$ ($1 \leq r \leq |T|$) to denote that, after an ID assignment f , the ID for path set t_j is r (or ranks the r -th among all the IDs). With f , we permute columns on \mathbf{M} and achieve \mathbf{N} . Column r in \mathbf{N} corresponds to column t_j in \mathbf{M} , i.e., $[n_{1r}, n_{2r}, \dots, n_{|S|r}]^T = [m_{1j}, m_{2j}, \dots, m_{|S|j}]^T$.

$$\mathbf{N} = f(\mathbf{M}) = \begin{matrix} & 1 & 2 & 3 & \dots & |T| \\ \begin{matrix} s_1 \\ s_2 \\ s_3 \\ \vdots \\ s_{|S|} \end{matrix} & \begin{pmatrix} n_{11} & n_{12} & n_{13} & \dots & n_{1|T|} \\ n_{21} & n_{22} & n_{23} & \dots & n_{2|T|} \\ n_{31} & n_{32} & n_{33} & \dots & n_{3|T|} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ n_{|S|1} & n_{|S|2} & n_{|S|3} & \dots & n_{|S||T|} \end{pmatrix} \end{matrix}$$

With matrix \mathbf{N} , we can easily calculate the number of \mathcal{AIB} s on each switch. For example, to compute it on switch s_i , we only need to sequentially check the elements on the i -th row of \mathbf{N} . If two sequential non-zero elements are the same, it indicates two consecutive IDs share the same egress port, and thus belong to the same \mathcal{AIB} . Otherwise, one more \mathcal{AIB} is needed. Therefore, the total number of \mathcal{AIB} s on switch s_i is:

$$\mathcal{AIB}(s_i) = 1 + \sum_{r=1}^{|T|-1} (n_{ir} \oplus n_{i(r+1)}) \quad (1)$$

where $u \oplus v = 1$ if $u \neq v$ (0 is not counted in Equation 1), and 0 otherwise. With Equation 1, we can get the largest number of \mathcal{AIB} s (i.e., maximal routing entries) across all switches:

$$\mathcal{MAIB} = \max_{1 \leq i \leq |S|} \{\mathcal{AIB}(s_i)\} \quad (2)$$

With above, we can mathematically define our problem as finding an ID assignment or ordering $f: \mathbf{M} \rightarrow \mathbf{N}$

so that \mathcal{MAIB} is minimized:

$$\text{Objective: } \textit{minimize } \mathcal{MAIB} \quad (3)$$

3.2.2 Practical solution

Our problem is NP-complete as it can be reduced from the 3-SAT problem [37]. Thus, we resort to heuristics.

ID assignment algorithm: Blindly permuting columns on \mathbf{M} is brute-force and time-consuming. Our practical solution is guided by the following thinking. Each switch s_i has its own local optimal assignment f_i . But these individual local optimal assignments may conflict with each other by assigning different IDs to the same path set, causing ID inconsistency. Then, to generate a global optimized assignment from the local optimal ones, we resolve ID inconsistency for each path set incrementally. In other words, we require that each step of ID inconsistency correction introduces minimal additional new routing entries. Based on this consideration, we introduce ID_Assignment(\cdot) in Algorithm 1. The main idea behind the algorithm is as follows:

- *First, we optimally assign path set IDs in each individual switch locally.* We simply achieve this for each switch by assigning the path sets having the same egress ports with consecutive IDs (lines 1–2).
- *Second, we correct the inconsistent IDs for each path set on different switches.* After the first step, the ID for a path set on different switches may be different. We resolve this ID inconsistency using a greedy method: we pick one ID out of all the inconsistent IDs for a path set and let others be consistent with it provided that this incurs minimal new routing entries (lines 3–10). More specifically, in lines 6–9, we try each of the inconsistent IDs, calculate its \mathcal{MAIB} , and finally pick the one with the minimal \mathcal{MAIB} . The algorithm terminates after we resolve the ID inconsistency for all path sets.

Time complexity: Though not optimal, we found our algorithm is effective in compressing the routing entries in our evaluation. However, the problem is high time cost since it works on a large matrix. We intentionally designed Algorithm 1 to be of low time complexity, i.e., $O(|S|^2|T|)$ for the $|S| \times |T|$ matrix \mathbf{M} . Even though, we find that when the network scales to several thousands, it cannot return a result over 24 hours (see Table 5). Worse, it is possible that $|S| \sim 10^{4-5}$ and $|T| \sim 10^6$ or more for large DCNs. In such cases, even a linear time algorithm can be slow, let alone any advanced algorithms.

Speedup with equivalence reduction: We take advantage of DCN topology characteristics to speed up the runtime of our algorithm to $O(c^2|T|)$, while maintaining comparable results as the original algorithm. Here, c is the number of equivalence classes of a topology.

Algorithm 1 Coordinated ID assignment algorithm

```
ID_Assignment(M) /* M is initial matrix, N is
output */;
1 foreach row  $i$  of M (i.e., switch  $s_i$ ) do
2   assign path sets  $t_j (1 \leq j \leq |T|)$  having the same
    $m_{ij}$  values (i.e., ports) with consecutive IDs;
/* path sets are optimally encoded in each local
switch, but causing ID inconsistency across
switches, i.e., one path set may have different IDs
on different switches */;
3 N  $\leftarrow$  M with IDs specified for each  $t_j$  in each  $s_i$ ;
4 foreach column  $j$  of N (i.e., path set  $t_j$ ) do
5   if  $t_j$  has inconsistent IDs then
6     let  $C = \{c_1, c_2, \dots, c_k\}, (1 < k \leq |S|)$  be
     the inconsistent IDs;
7     foreach  $c \in C$  do
8       use  $c$  to resolve the inconsistency by re-
       placing all other  $c_i$  in  $C$  with  $c$  on each
       relevant switch;
9       compute  $\mathcal{MAIB}$  using Equation 2;
10    set  $\text{ID}(t_j) \leftarrow c$  with the minimal  $\mathcal{MAIB}$ ;
11 return N /* N is inconsistency-free */;
```

DCN structures	Fattree	BCube	VL2	HyperX
$c, \#$ of equivalence classes	3	2	3	1
Time complexity	$O(T)$	$O(T)$	$O(T)$	$O(T)$

Table 3: Runtime of ID assignment algorithm with equivalence reduction on well-known DCNs.

The key observation is that most DCN topologies are regular and many nodes are equivalent (or symmetric). These equivalent nodes are likely to have similar numbers of routing states for a given ID assignment, especially when the paths are symmetrically distributed. For example, in the Fattree topology of Fig. 1, switches at the same layer are symmetric and are very likely to have similar number of routing entries. Thus, our hypothesis is that, by picking a representative node from each equivalence node class, we are able to optimize the routing tables for all nodes in the topology while spending much less time.

Based on the hypothesis, we significantly improve the runtime of Algorithm 1 with equivalence reduction. This speedup makes no change to the basic procedure of Algorithm 1. Instead of directly working on M with $|S|$ rows, the key idea is to derive a smaller M^* with c rows from M using equivalence reduction, and then apply $\text{ID_Assignment}(\cdot)$ on M^* . To this end, we first need to compute the equivalence classes among all the nodes, and there are many fast algorithms available for this purpose [14, 27]. This improvement enables our algorithm to finish computation within $O(|T|)$ time for various

well-known DCNs as shown in Table 3.

Our evaluation below validates our hypothesis and show that our method of equivalence reduction maintains good performance while consuming much less time.

3.3 Scalability evaluation

Evaluation setting: We evaluate XPath’s scalability on 4 well-known DCNs: Fattree [4], VL2 [16], BCube [17], and HyperX [3]. Among these DCNs, BCube is a server-centric structure where servers act not only as end hosts but also relay nodes for each other. For the other 3 DCNs, switches are the only relay nodes and servers are connected to ToRs at last hop. For this reason, we consider the paths between servers in BCube and between ToRs in Fattree, VL2 and HyperX.

For each DCN, we vary the network size (Table 4). We consider $k^2/4$ paths between any two ToRs in Fattree(k), $(k+1)$ paths between any two servers in BCube(n, k), D_A paths between any two ToRs in VL2(D_A, D_I, T), and L paths between any two ToRs in HyperX(L, S, T)⁵. These paths do not enumerate all possible paths in the topology, however, they cover all desired paths sufficient to exploit topology redundancy in each DCN.

Our scalability experiments run on a Windows server with an Intel Xeon E7-4850 2.00GHz CPU and 256GB memory.

Main results: Table 4 shows the results of XPath algorithm on the 4 well-known DCNs, which demonstrates XPath’s high scalability. Here, for paths to path sets aggregation we used CPF.

We find that XPath can effectively pre-install up to tens of billions of paths using tens of thousands of routing entries for very large DCNs. Specifically, for Fattree(64) we express 4 billion paths with 64K entries; for BCube(8,4) we express 5 billion paths with 47K entries; for VL2(100,96,100) we express 575 million paths with 117K entries; for HyperX(4,16,100) we express 17 billion paths with 36K entries. These results suggest that XPath can easily pre-install all desired paths into IP tables with 144K entries, and in the meanwhile XPath is still able to accommodate more paths before reaching 144K.

Understanding ID assignment: The most difficult part of the XPath compression algorithm is Step II (i.e., ID assignment), which eventually determines if XPath can pre-install all desired paths using 144K entries. The last two columns of Table 4 contrast the maximal entries be-

⁵Different DCNs use different parameters to describe their topologies. In Fattree(k), k is the number of switch ports; in BCube(n, k), n is the number of switch ports and k is the BCube layers; in VL2(D_A, D_I, T), D_A/D_I are the numbers of aggregation/core switch ports and T is the number of servers per rack; in HyperX(L, S, T), L is the number of dimensions, S is the number of switches per dimension, and T is the number of servers per rack.

DCNs	Nodes #	Links #	Original paths#	Max. entries # without compression	Path sets # after Step I compression	Max. entries # after Step I compression	Max. entries # after Step II compression
Fattree(8)	208	384	15,872	944	512	496	116
Fattree(16)	1,344	3,072	1,040,384	15,808	8,192	8,128	968
Fattree(32)	9,472	24,576	66,977,792	257,792	131,072	130,816	7,952
Fattree(64)	70,656	196,608	4,292,870,144	4,160,512	2,097,152	2,096,128	(<144K) 64,544
BCube(4, 2)	112	192	12,096	576	192	189	108
BCube(8, 2)	704	1,536	784,896	10,752	1,536	1,533	522
BCube(8, 3)	6,144	16,384	67,092,480	114,688	16,384	16,380	4,989
BCube(8, 4)	53,248	163,840	5,368,545,280	1,146,880	163,840	163,835	(<144K) 47,731
VL2(20, 8, 40)	1,658	1,760	31,200	6,900	800	780	310
VL2(40, 16, 60)	9,796	10,240	1,017,600	119,600	6,400	6,360	2,820
VL2(80, 64, 80)	103,784	107,520	130,969,600	4,030,400	102,400	102,320	49,640
VL2(100, 96, 100)	242,546	249,600	575,760,000	7,872,500	240,000	239,900	(<144K) 117,550
HyperX(3, 4, 40)	2,624	2,848	12,096	432	192	189	103
HyperX(3, 8, 60)	31,232	36,096	784,896	4,032	1,536	1,533	447
HyperX(4, 10, 80)	810,000	980,000	399,960,000	144,000	40,000	39,996	8,732
HyperX(4, 16, 100)	6,619,136	8,519,680	17,179,607,040	983,040	262,144	262,140	(<144K) 36,164

Table 4: Results of XPath on the 4 well-known DCNs.

fore and after our coordinated ID assignment for each DCN.

We find that XPath’s ID assignment algorithm can efficiently compress the routing entries by $2\times$ to $32\times$ for different DCNs. For example, before our coordinated ID assignment, there are over 2 million routing entries in the bottleneck switch (*i.e.*, the switch with the largest routing entries) for Fattree(64), and after it, we achieve 64K entries via prefix aggregation. In the worst case, we still compress the routing states from 240K to 117K in VL2(100,96,100). Furthermore, we note that the routing states may be further compressed using advanced IP routing table techniques such as [15, 34]. Our contribution is that the proposed coordinated ID assignment algorithm makes the prefix aggregation more efficient.

We note that our algorithm has different compression effects on different DCNs. As to the 4 largest topologies, we achieve a compression ratio of $\frac{2,096,128}{64,544} = 32.48$ for Fattree(64), $\frac{262,140}{36,164} = 7.25$ for HyperX(4,16,100), $\frac{163,835}{47,731} = 3.43$ for BCube(8,4), and $\frac{239,900}{117,550} = 2.04$ for VL2(100,96,100) respectively. We believe one important decisive factor for the compression ratio is the density of the matrix M . According to Equation 1, the number of routing entries is related to the non-zero elements in M . The sparser the matrix, the more likely we achieve better results. For example, in Fattree(64), a typical path set traverses $\frac{1}{32}$ aggregation switches and $\frac{1}{1024}$ core switches, while in VL2(100,96,100), a typical path set traverses $\frac{1}{2}$ aggregation switches and $\frac{1}{50}$ core switches. This indicates that M_{Fattree} is much sparser than M_{VL2} , which leads to the effect that the compression on Fattree is better than that on VL2.

Time cost: In Table 5, we show that equivalence reduction significantly speeds up the runtime of the ID assignment algorithm. For example, without equivalence reduction, it cannot return an output within 24 hours when the network scales to a few thousands. With it, we ob-

DCNs	Time cost (Second)	
	No equivalence reduction	Equivalence reduction
Fattree(16)	8191.121000	0.078000
Fattree(32)	>24 hours	4.696000
Fattree(64)	>24 hours	311.909000
BCube(8, 2)	365.769000	0.046000
BCube(8, 3)	>24 hours	6.568000
BCube(8, 4)	>24 hours	684.895000
VL2(40, 16, 60)	227.438000	0.047000
VL2(80, 64, 80)	>24 hours	3.645000
VL2(100, 96, 100)	>24 hours	28.258000
HyperX(3, 4, 60)	0.281000	0.000000
HyperX(4, 10, 80)	>24 hours	10.117000
HyperX(4, 16, 100)	>24 hours	442.379000

Table 5: Time cost of ID assignment algorithm with and without equivalence reduction for the 4 DCNs.

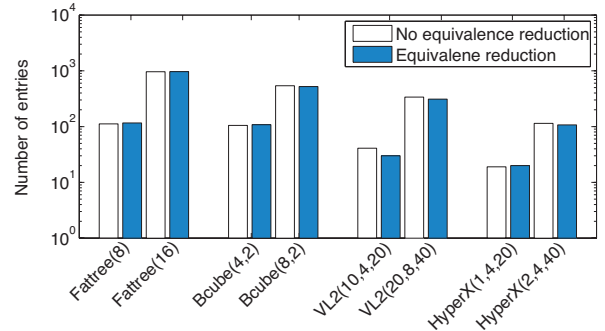


Figure 5: Effect of ID assignment algorithm with and without equivalence reduction for the 4 DCNs.

tain results for all the 4 DCNs within a few minutes even when the network is very large while maintaining good performance. Next, we show that equivalence reduction does not affect the compression effect of the algorithm.

Effect of equivalence reduction: In Fig. 5, we compare the performance of our ID assignment with and without equivalence reduction. With equivalence reduction, we use M^* (*i.e.*, part of M) to perform ID assignment, and it turns out that the achieved entries# is similar with that of without equivalence reduction (*i.e.*, use original M). This demonstrates the effectiveness of equivalence

reduction, which also validates our hypothesis in §3.2.2.

Results on randomized DCNs: We note that most other DCNs such as CamCube [2] and CiscoDCN [13] are regular and XPath can perform as efficiently as above. In recent work such as Jellyfish [39] and SWDC [38], people also discussed random graphs for DCN topologies. XPath’s performance is indeed unpredictable for random graphs. But for all the Jellyfish topologies we tested, in the worst case, XPath still manages to compress over 1.8 billion paths with less than 120K entries. The runtime varies from tens of minutes to hours or more depending on the degree of symmetry of the random graph.

4 Implementation and Experiments

We have implemented XPath on both Windows and Linux platforms, and deployed it on a 54-server Fattree testbed with commodity off-the-shelf switches for experiments. This paper describes the implementation on Windows. In what follows, we first introduce path ID resolution (§4.1) and failure handling (§4.2). Then, we present testbed setup and basic XPath experiments (§4.3).

4.1 Path ID resolution

As introduced in §2.2, path ID resolution addresses how to resolve the path IDs (*i.e.*, routing IPs) to a destination. To achieve fault-tolerant path ID resolution, there are two issues to consider. First, how to distribute the path IDs of a destination to the source. The live paths to the destination may change, for example, due to link failures. Second, how to choose the path for a destination, and enforce such path selection in existing networks.

These two issues look similar to the name resolution in existing DNS. In practice, it is possible to return multiple IPs for a server, and balance the load by returning different IPs to the queries. However, integrating the path ID resolution of XPath into existing DNS may challenge the usage of IPs, as legacy applications (on socket communication) may use IPs to differentiate the servers instead of routing to them. Thus, in this paper, we develop a clean-slate XPath implementation on the XPath manager and end servers. Each server has its original name and IP address, and the routing IPs for path IDs are not related to DNS.

To enable path ID resolution, we implemented a XPath software module on the end server, and a module on the XPath manager. The end server XPath software queries the XPath manager to obtain the updated path IDs for a destination. The XPath manager returns the path IDs by indexing the IP-to-ID mapping table. From the path IDs in the query response, the source selects one for the current flow, and caches all (with a timeout) for subsequent communications.

To maintain the connectivity to legacy TCP/IP stacks, we design an IP-in-IP tunnel based implementation.

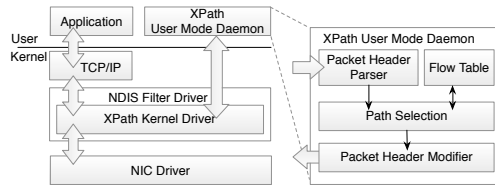


Figure 6: The software stacks of XPath on end servers.

The XPath software encapsulates the original IP packets within an IP tunnel: the path ID is used for the tunnel IP header and the original IP header is the inner one. After the tunnel packets are decapsulated, the inner IP packets are delivered to destinations so that multi-path routing by XPath is transparent to applications. Since path IDs in Fattree end at the last hop ToR, the decapsulation is performed there. The XPath software may switch tunnel IP header to change the paths in case of failures, while for applications the connection is not affected. Such IP-in-IP encapsulation also eases VM migration as VM can keep the original IP during migration.

We note that if VXLAN [42] or NVGRE [31] is introduced for tenant network virtualization, XPath IP header needs to be the outer IP header and we will need 3 IP headers which looks awkward. In the future, we may consider more efficient and consolidated packet format. For example, we may put path ID in the outer NVGRE IP header and the physical IP in NVGRE GRE Key field. Once the packet reaches the destination, the host OS then switches the physical IP and path ID.

In our implementation, the XPath software on end servers consists of two parts: a Windows Network Driver Interface Specification (NDIS) filter driver in kernel space and a XPath daemon in user space. The software stacks of XPath are shown in Fig. 6. The XPath filter driver is between the TCP/IP and the Network Interface Card (NIC) driver. We use the Windows filter driver to parse the incoming/outgoing packets, and to intercept the packets that XPath is interested in. The XPath user mode daemon is responsible to perform path selection and packet header modification. The function of the XPath filter driver is relatively fixed, while the algorithm module in the user space daemon simplifies debugging and future extensions.

From Fig. 6, we observe that the packets are transferred between the kernel and user space, which may degrade the performance. Therefore, we allocate a shared memory pool by the XPath driver. With this pool, the packets are not copied and both the driver and the daemon operate on the same shared buffer. We tested our XPath implementation (with tunnel) and did not observe any visible impact on the TCP throughput at the Gigabit line rate.

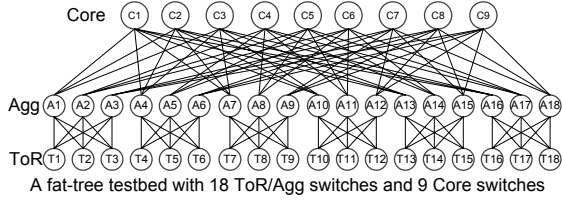


Figure 7: Fattree(6) testbed with 54 servers. Each ToR switch connects 3 servers (not drawn).

4.2 Failure handling

As introduced in §2.2, when a link fails, the switches on the failed link will notify the XPath manager. In our implementation, the communication channel for such notification is out-of-band. Such out-of-band control network and the controller are available in existing production DCNs [44].

The path IDs for a destination server are distributed using a query-response based model. After the XPath manager obtains the updated link status, it may remove the affected paths or add the recovered paths, and respond to any later query with the updated paths.

We implemented a failure detection method for TCP connection on the servers. In our XPath daemon, we check the TCP sequence numbers and switch the path ID once we detect that the TCP has retransmitted a data packet after a TCP timeout. The motivation is that the TCP connection is experiencing bad performance on the current path (either failed or seriously congested) and the XPath driver has other alternative paths ready for use. We note that there are faster failure detection mechanisms such as BFP or F10 [26] that can detect failures with $30\mu\text{s}$ (which we do not implement in our testbed), and XPath can be combined with them to perform fast rerouting. The key is that XPath does not require route re-convergence and is loop-free during failure recovery. This is because XPath pre-installs the backup paths and there is no need to do table re-computation except in exceptional cases where all backup paths are down.

4.3 Testbed setup and basic experiments

Testbed setup: We built a testbed with 54 servers connected by a Fattree(6) network (as shown in Fig. 7) using commodity Pronto Broadcom 48-port Gigabit Ethernet switches. On the testbed, there are 18 ToR, 18 Agg, and 9 Core switches. Each switch has 6 GigE ports. We achieve these 45 virtual 6-port GigE switches by partitioning the physical switches. Each ToR connects 3 servers; and the OS of each server is Windows Server 2008 R2 Enterprise 64-bit version. We deployed XPath on this testbed for our experimentation.

IP table configuration: In our testbed, we consider 2754 explicit paths between ToRs (25758 paths between end hosts). After running the two-step compression algorithm, the number of routing entries for the switch IP

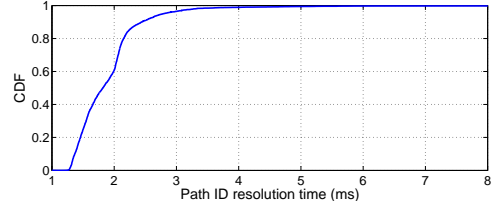


Figure 8: The CDF of path ID resolution time.

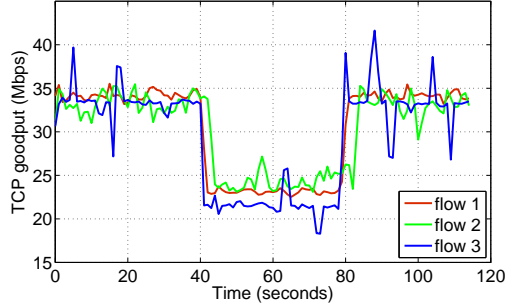


Figure 9: TCP goodput of three connections versus time on three phases: no failure, in failure, and recovered.

tables are as follows, ToR: 31~33, Agg: 48, and Core: 6. Note that the Fattree topology is a symmetric graph, the number of routing entries after compression is almost the same for switches at the same layer, which confirms our hypothesis in §3.2.2 that equivalent nodes are likely to have a similar number of entries.

Path ID resolution time: We measure the path ID resolution time at the XPath daemon on end servers: from the time when the query message is generated to the time the response from the XPath manager is received. We repeat the experiment 4000 times and depict the CDF in Fig. 8. We observe that the 99-th percentile latency is 4ms. The path ID resolution is performed for the first packet to a destination server that is not found in the cache, or cache timeout. A further optimization is to perform path ID resolution in parallel with DNS queries.

XPath routing with and without failure: In this experiment, we show basic routing of XPath, with and without link failures. We establish 90 TCP connections from the 3 servers under ToR T1 to the 45 servers under ToRs T4 to T18. Each source server initiates 30 TCP connections in parallel, and each destination server hosts two TCP connections. The total link capacity from T1 is $3 \times 1\text{G} = 3\text{G}$, shared by 90 TCP connections.

Given the 90 TCP connections randomly share 3 up links from T1, the load should be balanced overall. At around 40 seconds, we disconnect one link (T1 to A1). We use TCP sequence based method developed in §4.2 for automatic failure detection and recovery in this experiment. We then resume the link at time around 80 seconds to check whether the load is still balanced. We log the goodput (observed by the application) and show

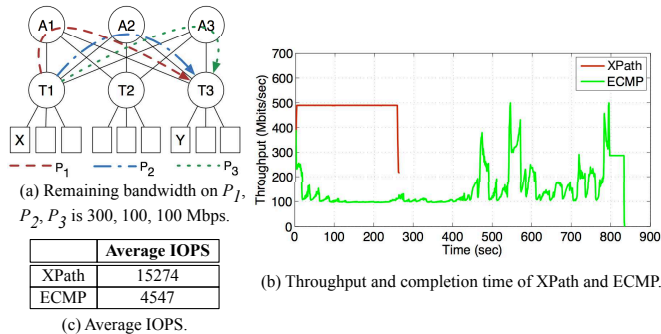


Figure 10: XPath utility case #1: we leverage XPath to make necessary bandwidth easier to implement for provisioned IOPS.

the results for three connections versus time in Fig. 9. Since we find that the throughput of all 90 TCP connections are very similar, we just show the throughput of one TCP connection for each source server.

We observe that all the TCP connections can share the links fairly with and without failure. When the link fails, the TCP connections traversing the failed link (T1 to A1) quickly migrate to the healthy links (T1 to A2 and A3). When the failed link recovers, it can be reused on a new path ID resolution after the timeout of the local cache. In our experiment, we set the cache timeout value as 1 second. However, one can change this parameter to achieve satisfactory recovery time for resumed links. We also run experiments for other traffic patterns, *e.g.*, ToR-to-ToR and All-to-ToR, and link failures at different locations, and find that XPath works as expected in all cases.

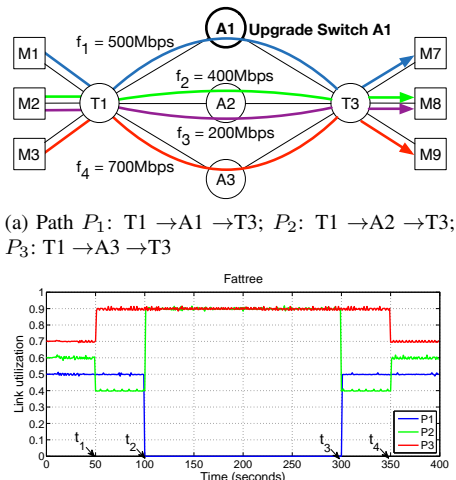
5 XPath Applications

To showcase XPath’s utility, we use it for explicit path support in 5 applications. The key is that, built on XPath, applications freely choose which path to use without worrying about how to setup the path and the time cost or overhead of setting up the path. In this regard, XPath emerges as an interface for applications to use explicit paths conveniently but not makes any choice on their behalf.

5.1 XPath for provisioned IOPS

In cloud services, there is an increasing need for provisioned IOPS. For example, Amazon EBS enforces provisioned IOPS for instances to ensure that disk resources can be accessed with high and consistent I/O performance whenever you need them [24]. To enforce such provisioned IOPS, it should first provide necessary bandwidth for the instances [10]. In this experiment, we show XPath can be easily leveraged to use the explicit path with necessary bandwidth.

As shown in Fig. 10(a), we use background UDP flows to stature the ToR-Agg links and leave the remaining bandwidth on 3 paths (P_1, P_2 and P_3) between X-Y as 300Mbps, 100Mbps, and 100Mbps respectively. Sup-



(a) Path P_1 : T1 \rightarrow A1 \rightarrow T3; P_2 : T1 \rightarrow A2 \rightarrow T3; P_3 : T1 \rightarrow A3 \rightarrow T3

(b) Time t_1 : move f_3 from P_2 to P_3 ; t_2 : move f_1 from P_1 to P_2 ; t_3 : move f_1 from P_2 to P_1 ; t_4 : move f_3 from P_3 to P_2 .

Figure 11: XPath utility case #2: we leverage XPath to assist zUpdate [25] to accomplish DCN update with zero loss.

pose there is a request for provisioned IOPS that requires 500Mbps necessary bandwidth (The provisioned IOPS is about 15000 and the chunk size is 4KB.). We now leverage XPath and ECMP to write 15GB data (\approx 4 million chunks) through 30 flows from X to Y, and measure the achieved IOPS respectively. The storage we used for the experiment is Kingston V+200 120G SSD, and the I/O operations on the storage are sequential read and sequential write.

From Fig. 10(c), it can be found that using ECMP we cannot provide the necessary bandwidth between X-Y for the provisioned IOPS although the physical capacity is there. Therefore, the actual achieved IOPS are only 4547, and the write with ECMP takes much longer time than that with XPath as in Fig. 10(c). This is because ECMP performs random hashing and cannot specify the explicit path to use, thus it cannot accurately use the remaining bandwidth on each of the multiple paths for end-to-end bandwidth provisioning. However, XPath can be easily leveraged to enforce the required bandwidth by taking advantage of its explicit path control. With XPath, we explicitly control how to use the three paths and accurately provide 500Mbps necessary bandwidth, achieving 15274 IOPS.

5.2 XPath for network updating

In production datacenters, DCN update occurs frequently [25]. It can be triggered by the operators, applications and various networking failures. zUpdate [25] is an application that aims to perform congestion-free network-wide traffic migration during DCN updates with zero loss and zero human effort. In order to achieve its goal, zUpdate requires explicit routing path control over the underlying DCNs. In this experiment, we show how

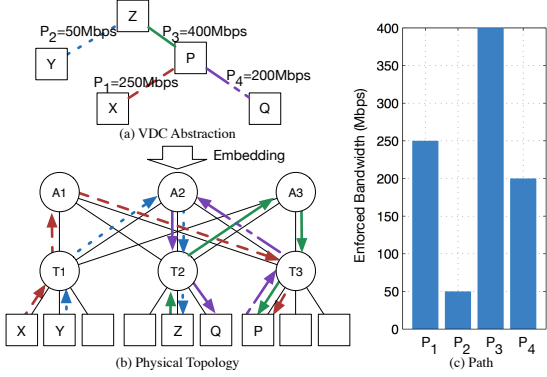


Figure 12: XPath utility case #3: we leverage XPath to accurately enforce VDC with bandwidth guarantees.

XPath assists zUpdate to accomplish DCN update and use a switch firmware upgrade example to show how traffic migration is conducted with XPath.

In Fig. 11(a), initially we assume 4 flows (f_1, f_2, f_3 and f_4) on three paths (P_1, P_2 and P_3). Then we move f_1 away from switch A_1 to do a firmware upgrade for switch A_1 . However, neither P_2 nor P_3 has enough spare bandwidth to accommodate f_1 at this point of time. Therefore we need to move f_3 from P_2 to P_3 in advance. Finally, after the completion of firmware upgrade, we move all the flows back to original paths. We leverage XPath to implement the whole movement.

In Fig. 11b, we depict link utilization dynamics of the whole process: at time t_1 , when f_3 is moved from P_2 to P_3 , the link utilization of P_2 drops from 0.6 to 0.4 and the link utilization of P_3 increases from 0.7 to 0.9. At time t_2 , when f_1 is moved from P_1 to P_2 , the link utilization of P_1 drops from 0.5 to 0 and the link utilization of P_2 increases from 0.4 to 0.9. The figure also shows the changes of the link utilization at time t_3 and t_4 when moving f_3 back to P_2 and f_1 back to P_1 . It is easy to see that with the help of XPath, P_1, P_2 and P_3 see no congestion and DCN update proceeds smoothly with 0 loss.

5.3 Virtual network enforcement with XPath

In cloud computing, virtual data center (VDC) abstraction with bandwidth guarantees is an appealing model due to its performance predictability in shared environments [7, 18, 45]. In this experiment, we show XPath can be applied to enforce virtual networks with bandwidth guarantees. We assume a simple SecondNet-based VDC model with 4 virtual links, and the bandwidth requirements on them are 50Mbps, 200Mbps, 250Mbps and 400Mbps respectively as shown in Fig. 12(a). We then leverage XPath’s explicit path control to embed this VDC into the physical topology.

In Fig. 12(b), we show that XPath can easily be employed to use the explicit paths in the physical topology

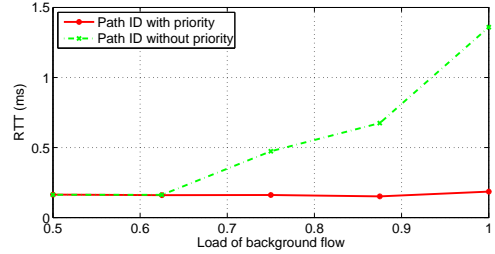


Figure 13: XPath utility case #4: we leverage XPath to express partition-aggregation traffic pattern (1-to- n, n -to-1) and enforce service differentiation.

with enough bandwidth to embed the virtual links. In Fig. 12(c), we measure the actual bandwidth for each virtual link and show that the desired bandwidth is accurately enforced. However, we found that ECMP cannot be used to accurately enable this because ECMP cannot control paths explicitly.

5.4 XPath support for partition-aggregation query

In web applications, the partition-aggregation paradigm is a foundation for many online services such as search query. They usually generate one-to-many and many-to-one communication patterns and has very demanding latency requirements. Using XPath, we can explicitly express such 1-to- n and n -to-1 patterns using $(n+1)$ path IDs, one ID for n -to-1 and n IDs for 1-to- n . These IDs form a group that can be leveraged for optimizations such as service differentiation.

In this experiment, we selected 9 servers in a pod to emulate a 1-to-8 (8-to-1) query-and-response, and we used 9 IDs to express such group communication patterns. We saturate the network with background traffic, and then leverage such 9 IDs to set priority to such query-and-response traffic.

In Fig. 13, we found that when the group of path IDs are referenced for priority, the query flows see persistently low RTTs of less than 200 μ s irrespective of the background traffic. However, if we do not set up a priority for these IDs, the RTT increases to the millisecond level as the background load increases. This experiment showcases XPath’s utility in service differentiation for partition-aggregation queries.

5.5 Map-reduce data shuffle with XPath

In Map-reduce applications, many-to-many data shuffle between the map and reduce stages can be time-consuming. For example, Hadoop traces from Facebook show that, on average, transferring data between successive stages accounts for 33% of the running times of jobs [12]. Using XPath, we can explicitly assign non-conflict parallel paths to speed up such many-to-many data shuffle. Usually, for a m -to- n data shuffle, we can use $(m+n)$ path IDs to express the communication patterns. The shuffle pattern can be well predicted using existing techniques [32].

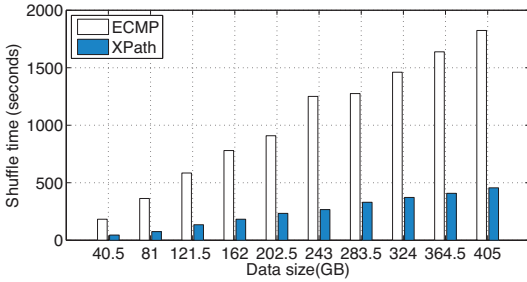


Figure 14: XPath utility case #5: we leverage XPath to select non-conflict paths to speed up many-to-many data shuffle.

In this experiment, we selected 18 servers in two pods of the Fattree to emulate a 9-to-9 data shuffle by letting 9 servers in one pod send data to 9 servers in the other pod. We varied the data volume from 40G to over 400G.

In Fig. 14, it can be seen that by using XPath for data shuffle, we can perform considerably better than randomized ECMP hash-based routing. More specifically, it reduces the shuffle time by over $3\times$ for most of the experiments. The reason is that XPath’s explicit path IDs can be easily leveraged to assign non-interfering paths for shuffling, thus the network bisection bandwidth is fully utilized for speedup.

6 Related Work

The key to XPath is its explicit path control. We note that many approaches such as source routing [36], MPLS [35]/ATM [9], OpenFlow [28] and the like, can also enable explicit path control. Each of them has its own limitation.

OpenFlow [28] has been used in many recent designs (e.g., [5, 8, 20, 21, 25]) to enable explicit path control. OpenFlow can establish fine-grained explicit routing path by installing flow entries in the switches via the OpenFlow controller. But in practice, there are dynamic flow setup overhead and small flow table entry number challenges to be solved. The on-chip OpenFlow forwarding rules in commodity switches are limited to a small number, typically 1–4K (2K in our Pronto switches). As a result, they need to dynamically identify and install small working sets of tunnels at different times [21, 25], which is technically complicated and heavy-weighted.

Source routing is usually implemented in software and slow paths, and not supported in the hardware of the data center switches, which typically only support destination IP based routing. Compared to source routing, XPath is readily deployable without waiting for new hardware capability; and XPath’s header length is fixed whereas it is variable for source routing with different path lengths.

In MPLS, paths (tunnels) can be explicitly setup before data transmission. Since the MPLS labels only have local significance, it requires a label distribution protocol

for label assignment. Further, MPLS is typically used only for traffic engineering in core networks instead of application-level or flow-level path control. In addition, it is known that the number of tunnels that existing MPLS routers can support is limited [6, 21].

SPAIN [29] builds a loop-free tree per VLAN and utilizes multiple paths across VLANs between two nodes, which increases the bisection bandwidth over the traditional Ethernet STP. However, SPAIN does not scale well because each host requires an Ethernet table entry per VLAN. Further, its network scale and path diversity are also restricted by the number of VLANs supported by Ethernet switches, e.g., 4096.

PAST [40] implements a per-address spanning tree routing for data center networks using the MAC table. PAST supports more spanning trees than SPAIN, but PAST does not support multi-paths between two servers, because a destination has only one tree. This is decided by MAC table size and its exact matching on flat MAC addresses.

Both SPAIN and PAST are L2 technologies. Relative to them, XPath builds on L3 and harnesses the fast-growing IP LPM table of commodity switches. One reason we choose IP instead of MAC is that it allows prefix aggregation. It is worth noting that our XPath framework contains both SPAIN and PAST. XPath can express SPAIN’s VLAN or PAST’s spanning tree using CPF, and it can also arrange paths using DPF and perform path ID encoding and prefix aggregation for scalability.

Finally, there are DCN routing designs that come with topologies, such as those introduced in Fattree [4], BCube [17], VL2 [16], PortLand [30], ALIAS [43], and so on. For example, the routings of Fattree and VL2 are only applicable to Clos topologies, ALIAS works on hierarchical topologies, and BCube and PortLand encode topology information into logical addresses (e.g., BCube IDs and PMAC) for routing. They are mostly customized to specific topologies and not easy to generalize. However, XPath works for arbitrary topologies.

7 Conclusion

XPath is motivated by the need for explicit path control in DCN applications. At its very core, XPath uses a path ID to express an end-to-end path, and pre-install all the desired path IDs between any s-d pairs into IP LPM tables using a two-step compression algorithm. Through extensive evaluation and implementation, we show that XPath is highly scalable and easy to implement in commodity switches. Finally, we used 5 primary experiments to show that XPath can directly benefit many popular DCN applications.

References

- [1] "Arista 7050QX," http://www.aristanetworks.com/media/system/pdf/Datasheets/7050QX-32_Datasheet.pdf, 2013.
- [2] H. Abu-Libdeh, P. Costa, A. Rowstron, G. O'Shea, and A. Donnelly, "Symbiotic Routing in Future Data Centers," in *SIGCOMM*, 2010.
- [3] J. H. Ahn, N. Binkert, A. Davis, M. McLaren, and R. S. Schreiber, "HyperX: Topology, Routing, and Packaging of Efficient Large-Scale Networks," in *SC*, 2009.
- [4] M. Al-Fares, A. Loukissas, and A. Vahdat, "A Scalable, Commodity Data Center Network Architecture," in *ACM SIGCOMM*, 2008.
- [5] M. Al-Fares, S. Radhakrishnan, B. Raghavan, N. Huang, and A. Vahdat, "Hedera: Dynamic Flow Scheduling for Data Center Networks," in *NSDI*, 2010.
- [6] D. Applegate and M. Thorup, "Load optimal MPLS routing with $N + M$ labels," in *INFOCOM*, 2003.
- [7] H. Ballani, P. Costa, T. Karagiannis, and A. Rowstron, "Towards Predictable Datacenter Networks," in *SIGCOMM*, 2011.
- [8] T. Benson, A. Anand, A. Akella, and M. Zhang, "MicroTE: Fine Grained Traffic Engineering for Data Centers," in *CoNEXT*, 2010.
- [9] U. Black, "ATM Volume I: Foundation for Broadband Networks," *Prentice Hall*, 1995.
- [10] I/O Characteristics. <http://docs.aws.amazon.com/AWSEC2/latest/UserGuide/ebs-io-characteristics.html>.
- [11] K. Chen, C. Guo, H. Wu, J. Yuan, Z. Feng, Y. Chen, S. Lu, and W. Wu, "Generic and Automatic Address Configuration for Data Centers," in *SIGCOMM*, 2010.
- [12] M. Chowdhury, M. Zaharia, J. Ma, M. Jordan, and I. Stoica, "Managing Data Transfers in Computer Clusters with Orchestra," in *SIGCOMM'11*.
- [13] Cisco, "Data center: Load balancing data center services," 2004.
- [14] P. T. Darga, K. A. Sakallah, and I. L. Markov, "Faster Symmetry Discovery using Sparsity of Symmetries," in *45st DAC*, 2008.
- [15] R. Draves, C. King, S. Venkatachary, and B. Zill, "Constructing optimal IP routing tables," in *INFOCOM*, 1999.
- [16] A. Greenberg, J. R. Hamilton, N. Jain, S. Kandula, C. Kim, P. Lahiri, D. A. Maltz, P. Patel, and S. Sengupta, "VL2: A Scalable and Flexible Data Center Network," in *ACM SIGCOMM*, 2009.
- [17] C. Guo, G. Lu, D. Li, H. Wu, X. Zhang, Y. Shi, C. Tian, Y. Zhang, and S. Lu, "BCube: A High Performance, Server-centric Network Architecture for Modular Data Centers," in *SIGCOMM*, 2009.
- [18] C. Guo, G. Lu, H. Wang, S. Yang, C. Kong, P. Sun, W. Wu, and Y. Zhang, "SecondNet: A Data Center Network Virtualization Architecture with Bandwidth Guarantees," in *CoNEXT*, 2010.
- [19] C. Guo, H. Wu, K. Tan, L. Shi, Y. Zhang, and S. Lu, "DCell: A scalable and fault-tolerant network structure for data centers," in *SIGCOMM'08*.
- [20] B. Heller, S. Seetharaman, P. Mahadevan, Y. Yiakoumis, P. Sharma, S. Banerjee, and N. McKeown, "ElasticTree: Saving Energy in Data Center Networks," in *NSDI*, 2010.
- [21] C.-Y. Hong, S. Kandula, R. Mahajan, M. Zhang, V. Gill, M. Nanduri, and R. Wattenhofer, "Achieving high utilization using software-driven WAN," in *ACM SIGCOMM*, 2013.
- [22] C. Hopps, "Analysis of an Equal-Cost Multi-Path Algorithm," *RFC 2992*, 2000.
- [23] Broadcom Strata XGS Trident II. <http://www.broadcom.com>.
- [24] Provisioned I/O-EBS. <https://aws.amazon.com/ebs/details>.
- [25] H. Liu, X. Wu, M. Zhang, L. Yuan, R. Wattenhofer, and D. Maltz, "zUpdate: Updating Data Center Networks With Zero Loss," in *ACM SIGCOMM*, 2013.
- [26] V. Liu, D. Halperin, A. Krishnamurthy, and T. Anderson, "F10: A Fault-Tolerant Engineered Network," in *NSDI*, 2013.
- [27] B. D. McKay, "Practical graph isomorphism," in *Congressus Numerantium*, 1981.
- [28] N. McKeown, T. Anderson, H. Balakrishnan, G. Parulkar, L. Peterson, J. Rexford, S. Shenker, and J. Turner, "OpenFlow: Enabling innovation in campus networks," *ACM CCR*, 2008.
- [29] J. Mudigonda, P. Yalagandula, and J. Mogul, "SPAIN: COTS Data-Center Ethernet for Multipathing over Arbitrary Topologies," in *NSDI*, 2010.
- [30] R. N. Mysore, A. Pamboris, N. Farrington, N. Huang, P. Miri, S. Radhakrishnan, V. Subramanya, and A. Vahdat, "PortLand: A Scalable Fault-Tolerant Layer 2 Data Center Network Fabric," in *SIGCOMM*, 2009.
- [31] NVGRE. <http://en.wikipedia.org/wiki/NVGRE>.
- [32] Y. Peng, K. Chen, G. Wang, W. Bai, Z. Ma, and L. Gu, "HadoopWatch: A First Step Towards Comprehensive Traffic Forecasting in Cloud Computing," in *INFOCOM*, 2014.
- [33] Announcing provisioned iops. <http://aws.amazon.com/about-aws/whats-new/2012/07/31/announcing-provisioned-iops-for-amazon-ebs/>.
- [34] G. Retvari, J. Tapolcai, A. Korosi, A. Majdan, and Z. Heszberger, "Compressing IP Forwarding Tables: Towards Entropy Bounds and Beyond," in *ACM SIGCOMM*, 2013.
- [35] E. Rosen, A. Viswanathan, and R. Callon, "Multiprotocol Label Switching Architecture," *RFC 3031*, 2001.
- [36] Source Routing. http://en.wikipedia.org/wiki/Source_routing.
- [37] Boolean satisfiability problem. http://en.wikipedia.org/wiki/Boolean_satisfiability_problem.

- [38] J.-Y. Shin, B. Wong, and E. G. Sirer, "Small-World Data-centers," in *ACM SoCC*, 2011.
- [39] A. Singla, C.-Y. Hong, L. Popa, and P. B. Godfrey, "Jellyfish: Networking Data Centers Randomly," in *NSDI*, 2012.
- [40] B. Stephens, A. Cox, W. Felter, C. Dixon, and J. Carter, "PAST: Scalable Ethernet for Data Centers," in *CoNEXT*, 2012.
- [41] Graph vertex coloring. http://en.wikipedia.org/wiki/Graph_coloring.
- [42] VXLAN. http://en.wikipedia.org/wiki/Virtual_Extensible_LAN.
- [43] M. Walraed-Sullivan, R. N. Mysore, M. Tewari, Y. Zhang, K. Marzullo, and A. Vahdat, "ALIAS: Scalable, Decentralized Label Assignment for Data Centers," in *SoCC*, 2011.
- [44] X. Wu, D. Turner, G. Chen, D. Maltz, X. Yang, L. Yuan, and M. Zhang, "NetPilot: Automating Datacenter Network Failure Mitigation," in *SIGCOMM*, 2012.
- [45] D. Xie, N. Ding, Y. C. Hu, and R. Kompella, "The Only Constant is Change: Incorporating Time-Varying Network Reservations in Data Centers," in *SIGCOMM*, 2012.