

# A Case for Performance- and Cost-aware Multi-Cloud Overlays

Bahador Yeganeh<sup>\*§</sup>, Ramakrishnan Durairajan<sup>†</sup>, Reza Rejaie<sup>†</sup> and Walter Willinger<sup>‡</sup>  
<sup>\*</sup> Snap, Inc.      <sup>†</sup> University of Oregon      <sup>‡</sup> NIKSUN, Inc.

**Abstract**—Modern enterprises are increasingly adopting multi-cloud strategies (i.e. connecting islands of resources from disparate public cloud providers, or CPs for short) due to benefits such as competitive pricing, global expansion, or improved reliability. While these benefits are compelling in their own right, what is critically lacking is a framework for establishing *optimized* multi-cloud overlays atop individual CP backbones in a performance- and cost-aware manner. A key challenge is that we have little understanding of the performance characteristics of CPs’ private backbones in light of multi-cloud overlays.

To address this challenge, we present a third-party, cloud-centric study to understand and examine the path, delay, and traffic-cost characteristics of CP backbones by deploying VMs in three global-scale CPs (i.e. AWS, Azure, and GCP). Our measurements reveal new insights including the “optimal backbone of cloud backbones” and a lack of path and delay asymmetries in it. Next, we report on several instances where performance-awareness of multi-cloud paths offer better latency reductions than default paths provided by CPs. While these results make a strong case for performance-aware multi-cloud overlays, the problem is further complicated by the varying transit costs/pricing models of CPs across different geographic regions. Based on our findings, we propose a research agenda for creating performance- and cost-aware multi-cloud overlays that deals with issues such as egress costs, considering IXPs as relays, using cloud auctions for transit cost pricing, and improving the performance of cloud-native applications.

**Index Terms**—Multi-cloud networks, Performance-aware Overlays, Network measurements

## I. INTRODUCTION

Modern enterprises are adopting multi-cloud strategies<sup>1</sup> at a rapid pace. Among the benefits of pursuing such strategies are competitive pricing, vendor lockout, global reach, and requirements for data sovereignty. According to a recent industry report, more than 85% of enterprises have already adopted multi-cloud strategies [1].

Despite this existing market push for multi-cloud strategies, we posit that there is a technology pull: *seamlessly connecting resources across disparate, already-competitive cloud providers (CPs) in a performance- and cost-aware manner is an open problem*. This problem is further complicated by two key issues. First, prior research on overlays has focused either on the public Internet [2] or on individual CP paths in isolation [3], [4], [5]. Second, because CP backbones are private and are invisible to traditional measurement techniques that

focus on public Internet paths, we lack a basic understanding of their performance, path, and traffic-cost characteristics.

This paper presents a cloud-centric measurement study that examines the performance and egress traffic costs of three prominent global-scale private cloud backbones (AWS, Azure, and GCP), from a third-party perspective. Our objective is to advocate for multi-cloud overlays that take into account path performance and egress costs of individual CP backbones.<sup>2</sup> Our measurements were run across 6 continents and 23 countries for 5 days (see Figure 1). Our measurements reveal several key insights. First, the cloud backbones (a) are optimal (i.e., 2x reduction in latency inflation ratio, which is defined as the ratio between line-of-sight and latency-based speed-of-light distances w.r.t. the public Internet), (b) lack path and delay asymmetry, and (c) are tightly interconnected with one another. Second, multi-cloud paths that are performance-aware exhibit higher latency reductions than default paths provided by CPs; e.g., 67% of all paths, 54% of all intra-CP paths, and 74% of all inter-CP paths experience an improvement in their latencies. Third, although traffic costs vary from location to location and across CPs, the costs are not prohibitively high. Based on these insights, we argue that enterprises and cloud users can indeed benefit from future efforts aimed at constructing high-performance overlay networks atop multi-cloud underlays in a performance- and cost-aware manner.

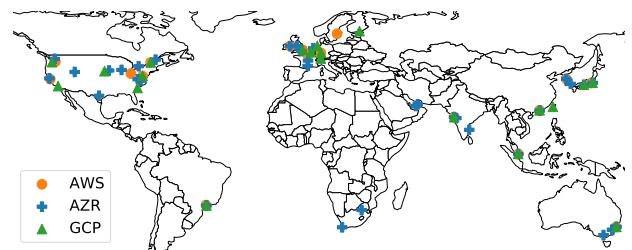


Fig. 1. Global regions for AWS, Azure, and GCP.

While our initial findings suggest that multi-cloud overlays are indeed beneficial for enterprises, establishing overlay-based connectivity to route enterprise traffic in a cost- and performance-aware manner among islands of disparate CP resources remains an open and challenging problem. The problem is further complicated by a lack of continuous multi-cloud

<sup>§</sup>Work done while a graduate student at the University of Oregon.

<sup>1</sup>This is different from hybrid cloud computing, where a direct connection exists between a public cloud and private on-premises enterprise server(s).

<sup>2</sup>While this study focuses on latency as the primary performance metric, we intend to broaden our investigation to include other performance metrics such as throughput and jitter in future work.

measurements [6], [7] and vendor-agnostic APIs [8]. In § VI, we outline our research agenda in pursuit of demonstrating the full potential of adopting multi-cloud strategies and discuss some of the open problems that arise in this context.

## II. CLOUD-CENTRIC MEASUREMENTS

In this section, we describe the CPs considered and the tools and datasets used in this study.

### A. Cloud Providers

In this study, we limit our focus to the top-3 cloud providers, namely Amazon Web Services (AWS), Microsoft Azure, and Google Cloud Platform (GCP) as they collectively account for the majority of the public cloud market share [9], [10]. We start by creating small VM instances within all regions of each CP. This resulted in a total of 63 regions (17, 26, and 20 regions for AWS, Azure, and GCP, respectively). We note that some regions are dedicated to government agencies and are not available to the public. Furthermore, we were not able to allocate VMs in 6 Azure regions<sup>3</sup> as those regions were either overpopulated or did not have free resources available at the time of this study. Next, we identify the physical locations of each CPs’ datacenters all over the world. We use a wide variety of data sources that point to exact (e.g., [11], [12], [13]) or approximate locations (e.g., [14], [15], [16], [17]) for these datacenters. We compile the information from these resources and, in the absence of information for a CP’s region, we default to the closest metro area which is reported by the CP. In summary, the VM instances used in our study are present in 6 continents, 23 countries, and in the proximity of 41 major metros as shown in Figure 1.

### B. Tools and Datasets Used

From each VM, we emit *paris-traceroute* and *ping* measurements using *scamper* [18] towards all other VM instances in successive 1-minute rounds for a duration of 5 days [19]. Subsequently, the resultant hops of each traceroute measurement are annotated with their corresponding ASN and ORG identifiers using prefix origins from the Routeviews [20] and the RIPE RIS [21]. The prefix origins are collected through BGPStream [22] and CAIDA’s AS-to-ORG dataset [23]. Furthermore, the existence of IXP hops along the path is checked by matching hop addresses against the set of IXP prefixes published by PeeringDB [24], Packet Clearing House (PCH) [25], and Hurricane Electric (HE) using CAIDA’s aggregate IXP dataset [26].

## III. ARE CLOUD BACKBONES OPTIMAL?

In this section, we look into the path, performance, and latency characteristics of CP backbones using the measurement setup described above.

<sup>3</sup>Central India, Canada East, France South, South Africa West, Australia Central, and Australia Central 2

### A. Path Characteristics of CP Backbones

As mentioned above, we measure the AS and ORG path for all of the collected traceroutes. In all our measurements, we observe multiple ASes for AWS *only* (AS14618 and AS16509). Hence, without loss of generality, from this point onward we only present statistics using the ORG measure. We measure the ORG-hop length for all unique paths and find that for 96% of our measurements, we only observe 2 ORGs (i.e. the source and destination CP networks). Out of the remaining paths, we observe that 3.85%, 0.03%, and 0.02% have 3, 4, and 5 ORG hops, respectively. These observations indicate two key results. First, all intra-CP measurements (and, hence, traffic) remain *almost always* within the CPs’ backbones. Second, the CP networks are tightly interconnected with each other and establish private peerings between each other on a global scale. Surprised by these findings, we take a closer look at the 4% of paths that include other networks along their paths. About 83% of these paths have a single IXP hop between the source and destination CPs. That is, the CPs are peering directly with each other over an IXP fabric. For the remaining 17% of paths, we observe two IXPs. Examples include paths (i) sourced from AWS in Seoul, KR, and destined to Azure in Johannesburg, ZA; and (ii) sourced from various GCP regions and destined to AWS in Hong Kong.

**Main findings:** *All intra-CP and the majority of inter-CP traffic remains within the CPs’ networks and is transmitted between the CPs’ networks over private and public peerings. CP’s backbones are tightly interconnected and can be leveraged for creating a global multi-cloud overlay.*

### B. Performance Characteristics of CP Backbones

Using the physical location of data centers for each CP, we measure the geo-distance between each pair of regions within a CP’s network using the Haversine distance [27] and approximate the optimal latency using the speed of light (SPL) constraints.<sup>4</sup> Figure 2 depicts the CDF of median latency inflation, which is defined as the ratio of median of the measured latency and SPL latency calculated using line-of-sight distances for each CP.

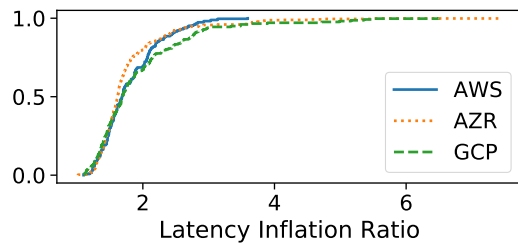


Fig. 2. Distribution of median latency inflation between network latency and RTT approximation using the speed of light constraints for all regions of each CP.

We observe median latency inflation of about 1.68, 1.61, and 1.68 for intra-CP paths of AWS, Azure, and GCP, respectively. Compared to a median latency inflation ratio of 3.2 for public Internet paths [28], these low latency inflation ratios attest to

<sup>4</sup>We use  $\frac{2}{3} * C$  for our calculations [28].

the optimal fiber paths and routes that are employed by CPs. Furthermore, Azure and GCP paths have long-tailed median latency inflation distributions while all intra-CP paths for AWS have a ratio of less than 3.6, making it the most optimal backbone among all CPs.

**Main findings:** *CPs employ an optimal fiber backbone with near line-of-sight latencies to create a global network. This result opens up a tantalizing opportunity to construct multi-cloud overlays in a performance-aware manner.*

### C. Latency Characteristics of CP Backbones

Next, we turn our attention to the latency characteristics of the CP backbones to create CP-specific latency profiles. Figure 3 shows the distribution of median RTT and standard deviation across different measurements for all paths between VM pairs. We observe a wide range of median RTT values between VM instances, which can be explained by the geographic distance between CP regions. Furthermore, latency between each pair is relatively stable across different measurements with a 90<sup>th</sup>-percentile standard deviation of less than 10ms—an observation consistent with Jain et al. [29].

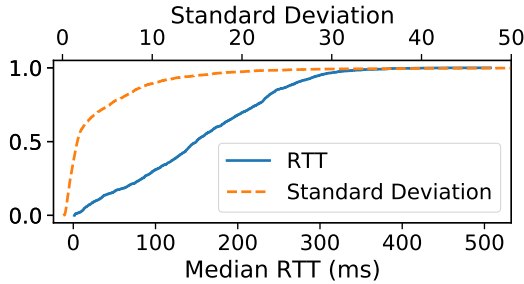


Fig. 3. Distribution of median RTT and standard deviation for latency measurements between all VM pairs.

In addition to stability characteristics, we also compare the forward and reverse path latencies by measuring the difference between the median of latencies in each direction. We find that paths exhibit symmetric latencies with a 95<sup>th</sup>-percentile latency difference of 0.06ms among all paths as shown in Figure 4.

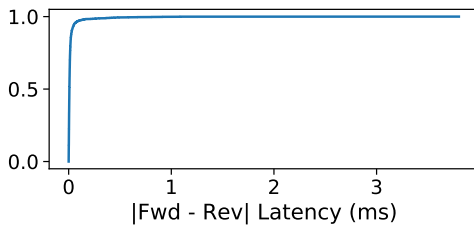


Fig. 4. Distribution for the difference in latency between forward and reverse directions for unique paths.

**Main findings:** *Cloud paths exhibit a stable and symmetric latency profile over our measurement period, making them ideal for reliable multi-cloud overlays.*

## IV. ARE THERE BENEFITS TO PERFORMANCE-AWARENESS IN MULTI-CLOUD PATH SELECTION?

While the default paths provided by individual CPs are performant in their own right, what are the performance

characteristics of paths that are indirect i.e. stitched together from two or more CPs (i.e. multi-cloud paths)? Are there opportunities to further improve the performance of multi-cloud paths by leveraging the performance characteristics of individual CP backbones? We seek to answer these questions in this section.

### A. Overall Latency Improvements

The distribution of latency reduction percentage for all, intra-CP, and inter-CP paths is shown in Figure 5. From this figure, we observe that about 67%, 54%, and 74% of all, intra-CP, and inter-CP paths experience an improvement in their latency using an indirect optimal path. These optimal paths can be constructed by relaying traffic through one or multiple intermediary CP regions. We provide more details on the intra- and inter-CP optimal overlay paths below.

To complement Figure 5, Figure 6-(left) shows the distribution of the number of relay hops along optimal paths. From this figure, we find that the majority (84%) of optimal paths can be constructed using *only* one relay hop while some paths can go through as many as 6 relay hops. Almost all of the optimal paths with latency reductions greater than 30% have less than 4 relay hops as shown in Figure 6-(right). In addition, we observe that the median latency reduction percentage increases with the number of relay hops. We note that (a) forwarding traffic through additional relay hops might have negative effects (e.g., increase in latencies) and (b) optimal paths with many relay hops might have an alternative path with fewer hops and comparable performance. We plan to study these two cases in more detail as part of our future work (§ VI).

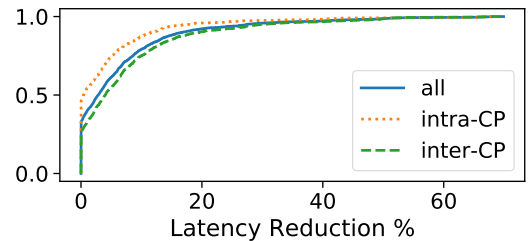


Fig. 5. Distribution for RTT reduction ratio through all, intra-CP, and inter-CP optimal paths.

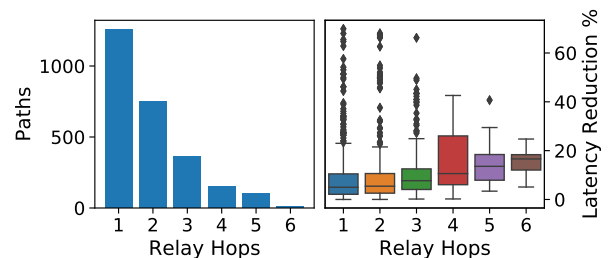


Fig. 6. Distribution for the number of relay hops along optimal paths (left) and the distribution of latency reduction percentage for optimal paths grouped based on the number of relay hops (right).

Lastly, we measure the prevalence of each CP along optimal paths and find that AWS, Azure, and GCP nodes are selected as relays for 97%, 52%, and 42% of optimal paths.

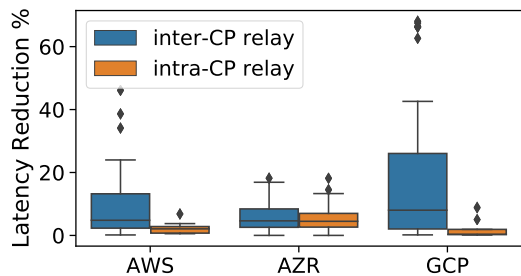


Fig. 7. Distribution of latency reduction percentage for intra-CP paths of each CP, divided based on the ownership of the relay node.

### B. Intra-CP Latency Improvements

We present statistics on the possibility of optimal overlay paths that are sourced and destined towards the same CP network (i.e. intra-CP overlays). Figure 7 depicts the distribution of latency reduction ratio for intra-CP paths of each CP. The distributions are grouped based on the owner (CP) of the relay node. From this figure, we see that intra-CP paths benefit more from relays located in other CP’s networks compared to relay nodes located in the same CP. More specifically, we find that 26%, 32%, and 12% of intra-CP paths for AWS, Azure, and GCP can benefit from relay nodes within their own network. This finding demonstrates that each CP can benefit from other CPs’ backbones to further improve their intra-network latency/performance.

**Main findings:** *Surprisingly, we find that intra-CP optimal paths can be constructed using relay hops that belong to a different CP.*

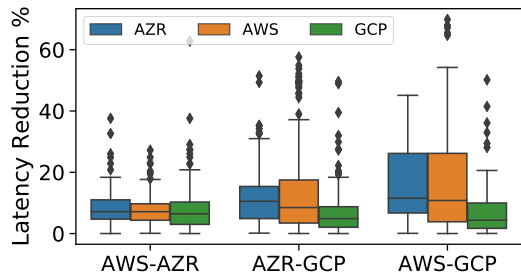


Fig. 8. Distribution of latency reduction ratio for inter-CP paths of each CP, divided based on the ownership of the relay nodes.

### C. Inter-CP Latency Improvements

We next focus on the possibility of overlay paths that are sourced from one CP and destined towards a different CP (i.e. inter-CP overlays). Figure 8 presents the latency reduction percentage for inter-CP paths. For brevity, only one direction of each CP pair is presented as the reverse direction is identical. We further divide each inter-CP path based on the relay nodes’ ownership. From this figure, we make several observations. First, optimal paths constructed using GCP nodes *as relays* exhibit the least amount of latency reduction. Second, *AWS-AZR* paths have lower values of latency reduction with equal amounts of reduction across each relay type. This is indicative of a tight coupling between these networks. Lastly,

optimal paths with AWS relays tend to have higher latency reductions which are in line with our observations in §III-A regarding AWS’ backbone.

**Main findings:** *Similar to intra-CP paths, inter-CP paths can benefit from relay nodes to construct new, optimal paths with lower latencies. Moreover, inter-CP paths tend to experience greater reductions in their latency.*

## V. CHALLENGES IN CREATING MULTI-CLOUD OVERLAYS

### A. Traffic Costs of CP Backbones

We turn our focus to the cost of sending traffic via CP backbones. Commonly, CPs charge their customers for traffic that is transmitted from their VM instances. That is, customers are charged *only* for egress traffic; all ingress traffic is free. Moreover, traffic is billed on a volume-by-volume basis (e.g., per GB of egress traffic) but each CP has a different set of rules and rates that govern their pricing policies. For example, we find that AWS and GCP have lower rates for traffic that remains within their network (i.e. is sourced and destined between different regions of their network) while Azure is agnostic to the destination of the traffic. Furthermore, GCP has different rates for traffic destined for the Internet based on the geographic region of the destination address. We compile all these pricing policies based on the information that each CP provides on their webpage [30], [31], [32] into a series of rules that allow us to infer the cost of transmitting traffic from each CP instance to other destinations.

**Traffic costs for AWS.** For AWS (see Figure 9), we observe that intra-CP traffic is always cheaper than inter-CP traffic except for traffic that is sourced from Australia and Korea. Furthermore, traffic sourced from the US, Canada, and European regions has the lowest rate while traffic sourced from Brazil has the highest charge rate per volume of traffic. Lastly, traffic is priced in multiple tiers defined based on the volume of exchanged traffic and we see that exchanging extra traffic leads to lower charging rates.

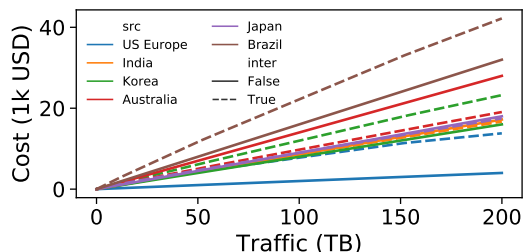


Fig. 9. Cost of transmitting traffic sourced from different groupings of AWS regions. Dashed (solid) lines present inter-CP (intra-CP) traffic cost.

**Traffic costs for Azure.** Azure’s pricing policy is more simple (see Figure 10). Global regions are split into multiple large-sized areas namely (i) North America and Europe excluding Germany, (ii) Asia and Pacific, (iii) South America, and (iv) Germany. Each of these areas has a different rate, with North America and Europe being the cheapest while traffic sourced from South America can cost up to 3x more than North America. Lastly, as mentioned earlier, Azure is agnostic

to the destination of traffic and does not differentiate between intra-CP and traffic destined to the Internet.

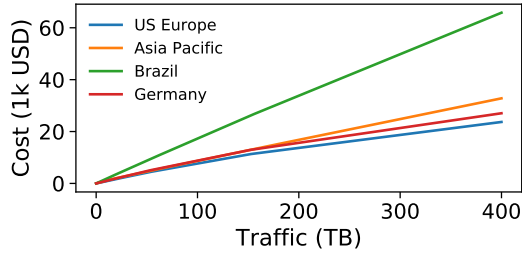


Fig. 10. Cost of transmitting traffic sourced from different groupings of Azure regions.

**Traffic costs for GCP.** GCP’s pricing policy is the most complicated among the top 3 CPs (see Figure 11). At a high level, GCP’s pricing policy can be determined based on (i) source region, (ii) destination geographic location, and (iii) whether the destination is within or outside GCP’s network or the Internet (intra-CP vs inter-CP). Intra-CP traffic generally has a lower rate compared to inter-CP traffic. Furthermore, traffic destined for China (excluding Hong Kong) and Australia has higher rates compared to other global destinations.

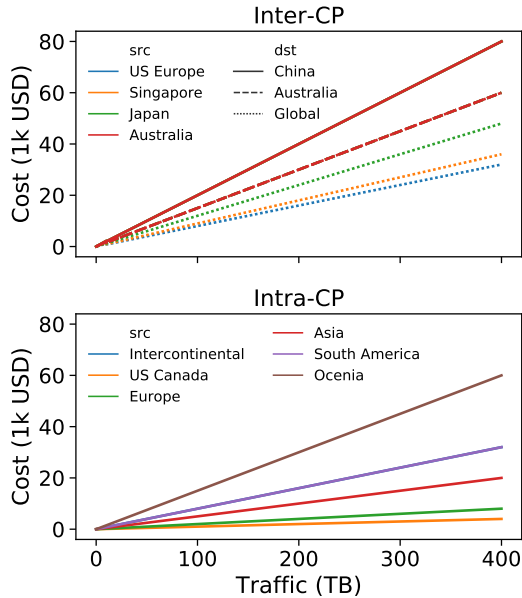


Fig. 11. Cost of transmitting traffic sourced from different groupings of GCP regions. Solid, dashed, and dotted lines represent the cost of traffic destined for China (excluding Hong Kong), Australia, and all other global regions accordingly.

### B. Cost Penalty for Multi-Cloud Overlays

Next, we seek an answer to the question of the cost incurred by using relay nodes from other CPs. Figure 12 depicts the distribution of cost penalty (i.e. the difference between the optimal overlay cost and default path cost) within various latency reduction percentage bins for transmitting 1TB of traffic.

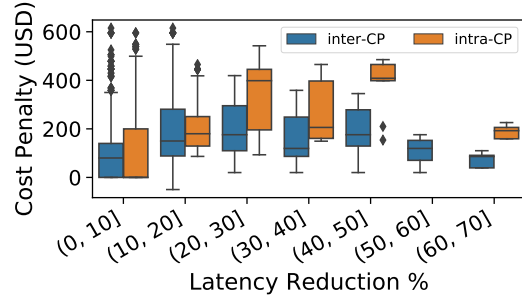


Fig. 12. Distribution of cost penalty within different latency reduction ratio bins for intra-CP and inter-CP paths.

From Figure 12, we make several key observations. First, we find that optimal paths between intra-CP endpoints incur higher cost penalties compared to inter-CP paths. This is expected as intra-CP paths tend to have lower charging rates and optimal overlays usually pass through a 3rd party CP’s backbone. Counter-intuitively, we next observe that the median cost penalty for paths with the most amount of latency reduction is less or equal to less optimal overlay paths. Lastly, we find that 3 of our optimal overlay paths have a negative cost penalty. That is, the optimal path costs are smaller than transmitting traffic directly between the endpoints. Upon closer inspection, we find that all of these paths are destined for the AWS Australia region and are sourced from GCP regions in Oregon US, Virginia US, and Montreal, Canada, respectively. All of these paths benefit from AWS’ lower transit cost toward Australia by handing off their traffic to a nearby AWS region. Motivated by this observation, for each set of endpoint pairs we find the path with the minimum cost. We find that the cost of traffic sourced from all GCP regions (except for GCP Australia) and destined to AWS Australia can be reduced by 28% by relying on AWS’ network as a relay hop. These cost-optimal paths on average experience a 72% inflation in their latency.

**Main findings:** *The added cost of overlay networks is not highly prohibitive. In addition to the inherent benefits of multi-cloud settings, our results demonstrate that enterprises and cloud users can construct high-performance overlay networks atop multi-cloud underlays in a cost-aware manner.*

## VI. RESEARCH AGENDA

Motivated by the above findings, we revisit the classical problem of creating overlay networks in a multi-cloud setting. To this end, we propose the following research agenda but leave the implementation details to future work.

### A. Constructing Multi-Cloud Overlays

**Performance- and Cost-aware Overlays.** Our first focus is to create an overlay network on top of multi-cloud underlays in a cost- and performance-aware manner. The starting point of our approach is to create a cloud-centric measurement service that continuously monitors the inter- and intra-CP links. Next, we plan to build vendor-agnostic APIs to connect the disparate island of CP resources. With the measurement service and APIs in place and given two locations (e.g., cities)

that are provided as input by a cloud user, the idea then is to construct a directed graph consisting of nodes that represent VM instances. Edges in the graph will be annotated with latencies and traffic-cost values from the measurement service.

Subsequently, we plan to build a framework that will find an optimal path in the directed graph that minimizes the latency of the path subject to a budget constraint for each pair of endpoints. This can be achieved by applying Dijkstra’s shortest path algorithm using the latency measures as the weight of each edge. Alternatively, we can formulate the problem as finding the *minimum cost flow* between a source (supply node) and destination (demand node) where links are annotated with their cost (latency) subject to an overall budget. We plan to consider this alternative as part of future work.

Supercloud [8] is the closest to our proposed work here. Supercloud proposes an abstraction layer to interconnect cloud providers and bridge the resources that are needed for storage and computation. Our proposed effort is complementary and goes beyond Supercloud by (i) using measurements to elucidate the performance issues of private multi-cloud underlays, and (ii) building overlays that are mindful of the performance and cost objectives in a multi-cloud setting.

**Robust Multi-Cloud Overlays.** The idea of creating multi-cloud overlays is both intriguing and promising. At the same time, we note that considering latency and cost *alone* can quickly turn relays into choke points and bottlenecks. In our ongoing analysis, we observed a skewed distribution regarding the use of CP regions along optimal paths (results are omitted due to space constraints). For example, about 25% of optimal paths use AWS India or Azure India as relays. Similarly, 15% of paths use AWS France as their relay hop. To tackle this problem, we plan to develop a *risk rank* for relay nodes. During overlay creation, the rank could be used as another input to the framework.

### B. Auctions for Multi-Cloud Overlays

To enable flexibility in creating multi-cloud overlays, we posit that CPs should adopt an auction-based model to lease their infrastructures—not just compute resources but also connectivity—to interested buyers/customers. Making the resources available via an auction recognizes the value of the multi-cloud overlay framework and measurement service (similar to the motivation for spot markets in cloud infrastructures). Similar to any market with competing entities, we hypothesize that the CPs compete based on factors including transit cost of relay nodes, geographical diversity, and robustness of established overlay paths.

**IXPs as relays to optimize overlay costs.** Motivated by the observation that of the 4% of paths with more than 2 ORG hops, 83% included a hop that is an IXP (as discussed in § III-A), we identify an opportunity to leverage IXPs [33] to save transit costs. This can be achieved in three steps. First, by targeting IXPs with more than 1 CP as a tenant and measuring the latency characteristics from each CP region to those IXPs. In the second step, the directed graph will be augmented with IXP nodes and their corresponding latency characteristics.

Third, the transit costs can be defrayed by accounting for the peering cost based on published peering costs by each CP [34], [35].

### C. Improving the Performance of Cloud-native Applications

While the first two items on our agenda focus on seamlessly creating overlays, many challenging issues remain. For one, comparing the performance of geo-distributed applications (e.g. Cassandra) on multi-cloud overlays (with cloud relay nodes) vs. the public Internet looms as an important open problem. Second, the recent trend of deploying microservices as a “single cluster” (e.g. using Docker Swarm) over overlay networks poses many challenges in a multi-cloud setting. Among these challenges are (a) how to create monitoring techniques in the face of CP-specific policies (e.g., some CPs might block measurement probes), (b) how to seamlessly adapt to the resource heterogeneity of underlays, and (c) how to create APIs to communicate the dynamism of underlays (e.g., failures, planned maintenances, etc.) to the cluster manager. We plan to address these problems as part of our future work.

## VII. SUMMARY

Market push indicates that the future of enterprises is multi-cloud. This market push has identified an existing “technology pull”—an acute need for a framework for seamlessly gluing the public cloud resources together in a cost- and performance-aware manner. A key reason behind this technology pull is that we lack a detailed understanding of the path, delay, and traffic-cost characteristics of the CPs’ private backbones. Our cloud-centric measurement study sheds light onto these characteristics and reveals several new/interesting insights into the CPs’ (private) backbones, including optimal cloud backbones, lack of delay and path asymmetries in cloud paths, possible latency improvements in inter- and intra-cloud paths, and traffic-cost characteristics. In short, this paper makes a strong case that the time is now for a “technology push” that realizes the full potential of multi-cloud strategies via overlays.

## ACKNOWLEDGEMENTS

We thank the anonymous reviewers for their insightful feedback. This work was supported by the National Science Foundation through CNS 2145813, CNS 1320977 and CNS 1719165. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of NSF.

## REFERENCES

- [1] A. Krishna, S. Cowley, S. Singh, and L. Kesterson-Townes, “Assembling your cloud orchestra: A field guide to multicloud management,” <https://www.ibm.com/thought-leadership/institute-business-value/report/multicloud>.
- [2] D. Andersen, H. Balakrishnan, F. Kaashoek, and R. Morris, “Resilient overlay networks,” in *Proceedings of ACM SOSP*, 2001.
- [3] O. Haq, M. Raja, and F. R. Dogar, “Measuring and improving the reliability of wide-area cloud paths,” in *proceedings of WWW*, 2017.
- [4] F. Lai, M. Chowdhury, and H. V. Madhyastha, “To relay or not to relay for inter-cloud transfers?” in *Workshop on Hot Topics in Cloud Computing*, 2018.

- [5] P. Costa, M. Migliavacca, P. Pietzuch, and A. L. Wolf, "Naas: Network-as-a-service in the cloud," in *Workshop on Hot Topics in Management of Internet, Cloud, and Enterprise Networks and Services*, 2012.
- [6] B. Yeganeh, R. Durairajan, R. Rejaie, and W. Willinger, "How cloud traffic goes hiding: A study of amazon's peering fabric," in *proceedings of IMC*, 2019.
- [7] —, "A first comparative characterization of multi-cloud connectivity in today's internet," in *proceedings of PAM*, 2020.
- [8] Z. Shen, Q. Jia, G.-E. Sela, W. Song, H. Weatherspoon, and R. Van Renesse, "Supercloud: A library cloud for exploiting cloud diversity," *ACM Transactions on Computer Systems (TOCS)*, vol. 35, no. 2, pp. 1–33, 2017.
- [9] Park my Cloud, "AWS vs Azure vs Google Cloud Market Share," <https://www.parkmycloud.com/blog/aws-vs-azure-vs-google-cloud-market-share/>, 2019.
- [10] Business Insider, "Goldman Sach forecasts the cloud computing market will keep growing," <https://bit.ly/2KG80G6>, 2019.
- [11] WikiLeaks, "Amazon Atlas," <https://wikileaks.org/amazon-atlas/>, 2018.
- [12] Build Azure, "Microsoft Azure Region Map," <https://map.buildazure.com/>, 2019.
- [13] Google, "Data center locations," <https://www.google.com/about/datacenters/inside/locations/index.html>, 2019.
- [14] R. Miller, "Regional Data Center Clusters Power Amazon's Cloud," <https://datacenterfrontier.com/regional-data-center-clusters-power-amazons-cloud/>, 2015.
- [15] I. Burrington, "Why Amazon's Data Centers Are Hidden in Spy Country," <https://www.theatlantic.com/technology/archive/2016/01/amazon-web-services-data-center/423147/>, 2016.
- [16] G. Plaven, "Amazon keeps building data centers in Umatilla, Morrow counties," <http://www.eastoregonian.com/eo/local-news/20170317/amazon-keeps-building-data-centers-in-umatilla-morrow-counties>, 2017.
- [17] M. Williams, "Amazon's central Ohio data centers now open," <http://www.dispatch.com/content/stories/business/2016/10/18/amazon-data-centers-in-central-ohio-now-open.html>, 2016.
- [18] M. Luckie, "Scamper: a scalable and extensible packet prober for active measurement of the internet," in *IMC*. ACM, 2010.
- [19] "Multi-cloud latencies dataset," [https://onrg.gitlab.io/projects/cloud\\_networking/#multi-cloud-network-management](https://onrg.gitlab.io/projects/cloud_networking/#multi-cloud-network-management).
- [20] University of Oregon, "Route Views Project," <http://www.routeviews.org/>.
- [21] RIPE, "RIPE RIS," 2019.
- [22] C. Orsini, A. King, D. Giordano, V. Giotsas, and A. Dainotti, "Bgp-stream: a software framework for live and historical bgp data analysis," in *proceedings of IMC*, 2016.
- [23] B. Huffaker, K. Keys, M. Fomenkov, and K. Claffy, "AS-to-Organization Dataset," <http://www.caida.org/research/topology/as2org/>, 2018.
- [24] PeeringDB, "PeeringDB," <https://www.peeringdb.com/>, 2019.
- [25] PCH, "Packet Clearing House," <https://www.pch.net/>, 2019.
- [26] CAIDA, "The CAIDA UCSD IXPs Dataset," 2019.
- [27] C. C. Robusto, "The cosine-haversine formula," *The American Mathematical Monthly*, 1957.
- [28] A. Singla, B. Chandrasekaran, P. Godfrey, and B. Maggs, "The internet at the speed of light," in *proceedings of HotNets*, 2014.
- [29] P. Jain, S. Kumar, S. Wooders, S. G. Patil, J. E. Gonzalez, and I. Stoica, "Skyplane: Optimizing transfer cost and throughput using cloud-aware overlays," in *proceedings of NSDI*, 2023.
- [30] Amazon, "EC2 Instance Pricing," <https://aws.amazon.com/ec2/pricing/on-demand/>, 2019.
- [31] Microsoft, "Bandwidth Pricing," <https://azure.microsoft.com/en-us/pricing/details/bandwidth/>, 2019.
- [32] Google, "Google Compute Engine Pricing," <https://cloud.google.com/compute/pricing#network>, 2019.
- [33] P. Sermpezis, G. Nomikos, and X. A. Dimitropoulos, "Re-mapping the Internet: Bring the IXPs into Play," *CoRR*, 2017.
- [34] Google, "Direct Peering — Interconnect," <https://cloud.google.com/interconnect/docs/how-to/direct-peering#pricing>, 2019.
- [35] Microsoft, "Virtual Network Pricing," <https://azure.microsoft.com/en-us/pricing/details/virtual-network/>, 2019.