

To Relay or Not to Relay for Inter-Cloud Transfers?

Fan Lai
University of Michigan

Mosharaf Chowdhury
University of Michigan

Harsha Madhyastha
University of Michigan

Abstract

Efficient big data analytics over the wide-area network (WAN) is becoming increasingly more popular. Current geo-distributed analytics (GDA) systems employ WAN-aware optimizations to tackle WAN heterogeneities. Although extensive measurements on public clouds suggest the potential for improving inter-datacenter data transfers via detours, we show that such optimizations are unlikely to work in practice. This is because the widely accepted mantra used in a large body of literature – WAN bandwidth has high variability – can be misleading. Instead, our measurements across 40 datacenters belonging to Amazon EC2, Microsoft Azure, and Google Cloud Platform show that the available WAN bandwidth is often spatially homogeneous and temporally stable between two virtual machines (VMs) in different datacenters, even though it can be heterogeneous at the TCP flow level. Moreover, there is little scope for either bandwidth or latency optimization in a cost-effective manner via relaying. We believe that these findings will motivate the community to rethink the design rationales of GDA systems and geo-distributed services.

1 Introduction

Many popular cloud applications nowadays host their services all around the world to meet the performance and regulatory requirements of their customers. Their customer-facing functionalities are often enabled by multiple internal services that depend on data and computation spread across the globe, which introduces new research issues at the intersections of networking, systems, and database. Consequently, a growing body of recent work has focused on enabling big data analytics across geo-distributed datacenters [1–4] – also known as geo-distributed analytics (GDA) – where computation is applied to in-place data at different sites. These GDA systems often involve large data transfers across the wide-area network (WAN) when they aggregate intermediate results generated at multiple sites.

Recent studies suggest that, in sharp contrast to intra-datacenter networks, inter-datacenter WANs experience high spatial and temporal bandwidth variations across different datacenters [1–7], and they impose new challenges in designing GDA systems. For example, they

suggest that the available WAN bandwidth between geographically-close regions can be up to $12\times$ more than that between distant regions [3]. To this end, state-of-the-art GDA systems propose heuristics to co-design intra-datacenter computation and inter-datacenter communication – subject to the scarce and heterogeneous WAN bandwidth – to improve GDA performance.

One possible implication of these WAN measurements between public cloud datacenters is that *tenants could potentially improve their WAN performance by relaying inter-datacenter transfers via detours*. Specifically, data transfers between two distant datacenters could be sped up by relaying the transfers via a third datacenter. For instance, WAN bandwidth measurements between cloud datacenters show that about 40% data transfers between Amazon EC2 datacenters can achieve more than $1.5\times$ bandwidth increase when they go through a one-hop relay instead of using the direct path. Such potential for relaying can seem more promising when it comes to multiple cloud providers, because the union of datacenters across multiple providers results in a geographically denser set of datacenters.

In this paper, we show that such optimizations are unlikely to work because the variations in WAN bandwidth shown in prior work are often artificial; as such, seemingly sensible rules of thumb – e.g., available bandwidth varies significantly across datacenter pairs – may be inappropriate. We argue that the measurement methodologies in prior work are misleading in at least two ways: (i) the revealed bandwidth heterogeneity is that of a single TCP connection – not the real bandwidth between two virtual machines (VMs) in different cloud datacenters – and this bandwidth is a function of the round-trip-time (RTT) between datacenters and the TCP window size; (ii) the available bandwidth between two VMs in different sites should be measured as the aggregate bandwidth of multiple TCP connections between them.

We substantiate our findings via extensive measurements on three popular cloud providers: Amazon EC2, Microsoft Azure, and Google Cloud Platform. Our experiments reconfirm some of the earlier findings (e.g., per-flow rate limiting), but provide some distinctive insights on the characteristics of inter-cloud WANs: (i) the available WAN bandwidth is homogeneous at the VM level and capped per-VM; (ii) the available WAN band-

System	Problem Statement
Clarinet [1]	Improve query response time via joint planning and scheduling, subject to heterogeneous bandwidth
Iridium [2]	Balance the transfer times among the heterogeneous WAN by optimizing data and task placement
Gaia [3]	Perform efficient machine learning across multiple datacenters, but WAN variations degrade performance
Amoeba [7]	Difficult to ensure timely data delivery, as WAN bandwidth varies in distance and time

Table 1: Selected recent GDA designs.

width is stable over periods of time; and (iii) bandwidth contention occurs at VMs instead of inter-datacenter links. Moreover, we show that there is not much scope for latency optimizations via relaying either. We believe that these findings will prove useful for practitioners, designers, and users of GDA systems.

Overall, we make the following contributions in this paper: (i) Extensive measurements on a multi-cloud deployment and novel insights regarding the characteristics of WAN across multiple cloud providers; (ii) An in-depth analysis of the opportunities to improve WAN performance by leveraging the power of multiple clouds, and we present several interesting angles for the design of GDA systems.

2 Related Work

Understanding the performance of WAN is crucial to several research communities. To a first approximation, our work can contribute to at least three categories.

Cloud Measurements Studies on bandwidth have shown prevalence and persistence of network variations over the WAN (e.g., [3, 5, 7]). These studies have either focused on the heterogeneity of bandwidth (e.g., [3, 6, 8]) or temporal WAN variations (e.g., [7]). Few studies provide in-depth analysis of inter-datacenter network across cloud providers (e.g., Azure and Google Cloud) and interpret those characteristics. Our work identifies the similarities and differences over three popular cloud providers.

WAN-Aware Optimization for GDA Recent work establishes the emerging problem of designing GDA systems (Table 1). Given WAN bandwidth heterogeneity, these works intelligently propose heuristics to design WAN-aware data analytics frameworks, showing very promising system-level performance improvements (e.g., [1, 2]) and WAN bandwidth usage reductions [9]. Throughout the paper, we refer to the insights provided by these studies and analyze the impacts of our new findings relevant to these GDA systems.

Overlay Routing on Clouds Geographically distributed datacenters connected by cloud providers’ backbone networks provide the opportunity for cloud users

Provider	VM Type	# of DCs Measured
Amazon	t2 micro	11
Microsoft	f1 micro	17
Google	n1-standard-1	12

Table 2: VM type and measured datacenters.

to construct performant overlay alternatives. Although extensive work have shown the effectiveness of overlay routing to protect against network outages or packet losses (e.g., [10–12]), they do not focus on improving network bandwidth or latency on clouds (e.g., VIA [13], ARROW [14]). Instead, we conduct a detailed analysis of the tenant-level opportunities to improve the WAN bandwidth and latency by relaying through multiple regions or clouds.

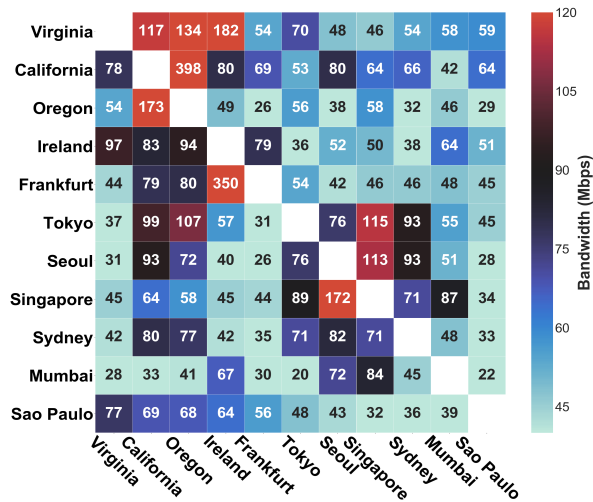
3 Inter-Cloud WAN Measurements

We start by presenting our measurement methodology. Then we quantify the bandwidth, latency, and opportunities for improving WAN performance via relaying.

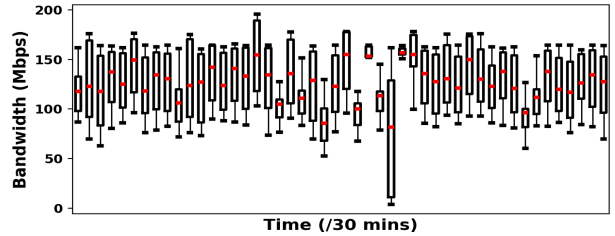
3.1 Measurement Methodology

Overview Our measurements span three popular cloud providers (Amazon, Microsoft, and Google) and aim to measure WAN bandwidth and latency of paths interconnecting VMs, including both intra-datacenter connections and inter-datacenter connections. This real cloud deployment can capture the performance that a tenant running a geo-distributed service in the cloud can expect. We conducted extensive measurements over 40 datacenter regions across three cloud providers’ datacenters and covered many trans-oceanic connections. Although our bandwidth measurements are limited by budget constraints, the results are generalized enough to support our findings. Details about the type of VMs and datacenters measured are noted in Table 2.

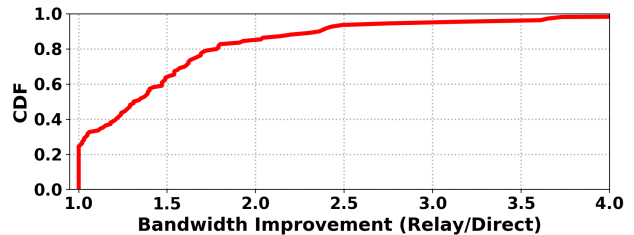
Measurement Tools We use the widely adopted measurement tool `iperf3` to measure the bandwidth between VMs. As already done in previous works [12], we use both ICMP-based `ping` and TCP-based `hping3` to measure the RTT between two VMs, because ICMP



(a) Bandwidth between 11 Amazon EC2 regions.



(b) Bandwidth from Virginia to California in a day.



(c) Throughput improvement via best relay.

Figure 1: WAN bandwidth on EC2 varies spatially and temporally when bandwidth is measured with a single TCP flow. In (b), the bottom and top ends of each boxplot represent the 25th and 75th percentiles, and the line in the middle represents the median. Samples are measured every 30 minutes, where each box reports the bandwidth variation over 60 seconds.

ping is not allowed on Azure. Note that results using ping and hping3 are similar on EC2 and Google Cloud.

3.2 Bandwidth

In this section, we first substantiate our findings by reproducing measurements from existing work. However, we show that observations in recent advances are misleading due to their measurement methodology, and we discuss possible reasons why WAN bandwidth was found to be heterogeneous and variable in prior measurements.

3.2.1 Prior Measurements are Misleading

We first measure the WAN bandwidth between 11 Amazon EC2 regions with the measurement methodology adopted in prior work. Specially, we use `iperf3` to measure the bandwidth of each pair of different regions every 15 minutes over a total period of 36 hours. Figure 1a reports the average WAN bandwidth between each pair of different regions, and Figure 1b shows the measured bandwidth from Virginia to California in 24 hours. Observations from these measurements are in line with recent work (e.g., [1, 3, 5–7]).

WAN Bandwidth Varies Greatly Across Different Pairs of Regions

As shown in Figure 1a, bi-directional bandwidth between the same pair of VMs is not strictly symmetric. For example, the average bandwidth from Oregon to Virginia is 146 Mbps, while the bandwidth of the opposite direction is 56 Mbps. Moreover, bandwidth between geographically-close regions (e.g., Frankfurt → Ireland) is up to 21× of the bandwidth between distant regions (e.g., Mumbai → Sao Paulo).

This bandwidth heterogeneity across the WAN has been accepted as a mantra in the literature.

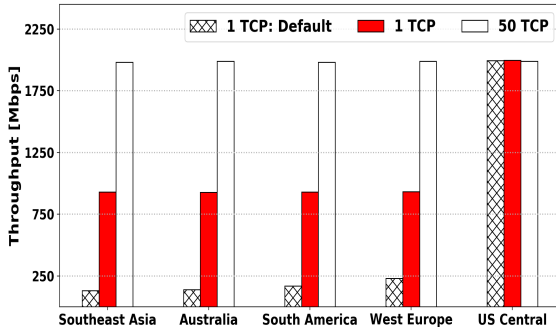
WAN Bandwidth Varies Temporally Measurement results in Figure 1b show that even bandwidth variations within every 60 seconds over 24 hours is highly dynamic, suggesting the difficulty to ensure expected data delivery in GDA systems. This intractable variation in inter-datacenter bandwidth has imposed a great challenge toward designing GDA systems [3].

Relaying Can Improve Throughput Observations from Figure 1c show that about 40% of data transfers on EC2 can obtain 1.5× or more bandwidth than the direct path when they go through a one-hop relay. This benefit becomes more promising (up to 10×) when tenants relay their traffic through the datacenters of other cloud providers.

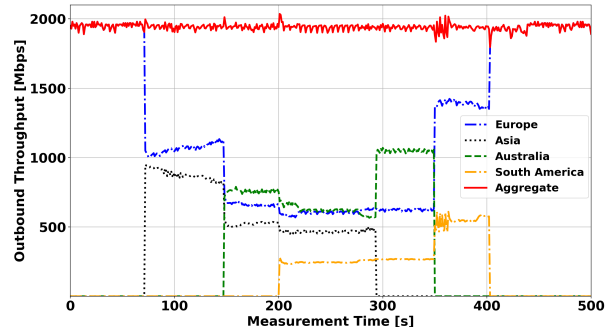
However, we argue that similar measurement results in the literature are misleading in two ways: (i) The variation and heterogeneity of WAN bandwidth are due to RTT variations and window sizes of the single TCP connection; (ii) The available bandwidth between two different sites should be considered as the aggregate bandwidth of multiple TCP connections. To generalize a step further, the well-known variability and heterogeneity of WAN bandwidth in prior work may be artificial! We substantiate these hypotheses next.

3.2.2 WAN Bandwidth is Often Homogeneous

In contrast to intra-datacenter networks, inter-datacenter WAN spreads across multiple regions in different conti-



(a) Throughput from other regions to US Center.



(b) Outbound throughput from US Center to other regions.

Figure 2: Available WAN bandwidth is homogeneous on Google Cloud. The default per-TCP throughput in (a) means the measured throughput without modifying the maximum window size in Linux. In (b), each line notes the aggregate throughput of 50 TCP connections. TCP connections from the VM in US Center to other sinks are created at different time.

nents. The RTT between two remote regions can reach up to 400 ms (see next section), which makes the WAN a long fat network with large bandwidth-delay products. In our measurements, we observe that the bandwidth-delay product of single TCP connections are nearly identical across different pairs of datacenters, confirming our hypothesis that the default TCP configuration in Linux is the culprit behind the misleading measurements.

In the rest of our measurements, we modify the TCP buffer size and window size in `/proc/sys/net/core` and `/proc/sys/net/ipv4`, thus achieving a maximum TCP window size of 7 MByte in `iperf3`. Furthermore, we study the available WAN bandwidth by setting up multiple TCP connections per VM. We conducted extensive experiments over multiple weeks on a subset of regions, and the results are quantitatively similar, without extra variations due to time of day impact. In the following, we provide a discussion of the obtained results contrasting them to existing literature.

Aggregate Outbound and Inbound Rates are Limited by the VM Rate Cap Most existing work for GDA treat the WAN as a full mesh and assume that bandwidth contentions occur at the uplinks or downlinks of VM pairs. However, our results from Figure 2b show that TCP connections of different inter-datacenter pairs will adaptively share the maximum available bandwidth when TCP connections are created or stopped, though their aggregate throughput is always around 2 Gbps. This means that the aggregate outbound rate of VMs is limited by the VM cap instead of the rate limiting associated with links. We have similar observations on the aggregate inbound rate of VMs. We hypothesize that all outbound traffics may travel through the same physical NIC.

Available Inter-Datacenter Bandwidth is Spatially Homogeneous Measurement results in Figure 2a reconfirm the rate-limiting on a per-flow basis in inter-

Type \ Region	Tokyo	Sydney	S.Paulo	Frankfurt	Oregon
t2 micro	413	420	432	426	441
m4 16xlarge	4740	4752	4761	4834	4787

Table 3: Throughput to US Center on EC2. (Mbps)

datacenter connections [7]. However, our results highlight three new aspects. First, although we notice a significant increase of per-TCP throughput on Azure and Google Cloud, such an increase does not show up on our measured EC2 VMs. We hypothesize that more conservative rate-limiting policies are used on EC2. Second, while the measured per-flow cap on EC2 is divergent among different region pairs, the throughput of single TCP connections across different regions can uniformly reach up to 1 Gbps on Azure and GCP. It suggests that changing TCP configurations does not work on our measured EC2 VMs. Third, if we focus on multiple TCP connections, the aggregate throughput is capped at the same limit for various inter-datacenter connections. As shown in Figure 2a and Table 3, although different VM types have divergent rate limiting, the maximum available WAN bandwidth from US Center to other remote regions is homogeneously capped per-VM.

Available Inter-Datacenter Bandwidth is Temporally Stable We further investigate the WAN bandwidth variations on these three cloud platforms. Figure 2b shows that the available bandwidth on Google Cloud is stable over 500 seconds. We note that this observation holds for Azure and EC2 high performance VMs in our sample measurements over multiple days, with no extra variations. This reconfirms our hypothesis that available inter-datacenter bandwidth is capped per-VM.

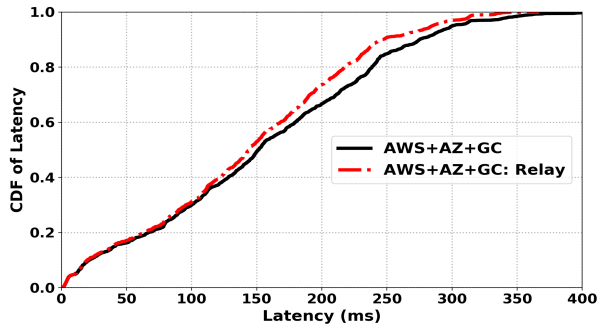


Figure 3: Latency across 40 datacenters.

Discussion The performance assessment presented above carries advantageous information to the design of WAN-aware cloud systems. In particular, the simplifying assumption in previous works that the runtime WAN environment of wide-area data analytics is temporally stable can be satisfied on some types of VMs, which makes it possible to ensure expected timely data delivery for data transfers. However, the VM rate cap on outbound traffics motivates us to focus on the bandwidth contentions inside VMs, instead of link-level optimizations used in existing GDA work.

3.3 Latency

We now present results of our measurement study of inter-datacenter latency. Our measurements of WAN latency span 40 datacenter regions across three cloud providers. It is worth noting that our measurements in August 2017 and January 2018 are quantitatively similar, as RTT mainly depends on the physical distance in the geo-distributed scenario. Figure 3 shows the measurement results.

WAN Latency Varies Greatly Between Different Regions Latencies between datacenters vary widely based on inter-datacenter distances. For example, latency between remote regions – e.g., Sao Paulo and Mumbai is up to 390 ms – imposes challenges for interactive cloud services across regions. Meanwhile, we observe that latency across the WAN is predictable with no significant variations over time (similar to [12]). The superior performance of inter-datacenter WAN over the public Internet is mainly due to the more manageable backbone networks on clouds.

4 To Relay or Not to Relay?

It is time to revisit the question: *Can we leverage the power of relay between clouds?* Prior work [12] has substantiated that detour routing is highly effective for cloud paths to improve network reliability. Unfortunately, there is little scope for higher available bandwidth and better

latency performance on the WAN by relaying.

Observations of the VM rate cap show that relay may not be helpful for higher network throughput. However, there are two interesting takeaways. First, it is encouraging for tenants to use multiple TCP connections for higher inter-datacenter throughput, because single TCP connections generally cannot saturate the per-VM rate cap. It is worth noting that multiple TCP connections will not result in extra cost for the same amount of inter-WAN data transfers, since cloud providers charge their users network traffic in terms of the volume of transfers. Second, when we consider EC2 applications that rely on single TCP connections, applications that can afford to pay extra can relay data transfers through Azure or Google Cloud to significantly improve network throughput. This is because the per-flow cap on our measured EC2 VMs is heterogeneous across different regions but greatly depends on latency, while the throughput of single TCP connections on Azure and Google Cloud are uniform around 1 Gbps.

On the other hand, there is little room to improve latency by relaying through different regions. Although the union of datacenters across multiple cloud providers results in a geographically denser set of datacenters than any single provider, datacenters on different clouds are geographically-close in the same continent. Hosting the service in a datacenter region that is close to users should be a better way to decrease network latency.

Finally, our measurements on clouds are far from comprehensive. Over 40 VM instance types with a variety of CPU, memory, disk, and network options are available on EC2 and Azure. Google provides 18 types and also allows customizing customer requirement for VMs memory and the number of CPU cores [15]. This flexible customization makes it pretty difficult to cover the entire VM-WAN performance landscape. In particular, the recently introduced “burstable” instances (e.g., t2.nano on EC2) that are significantly cheaper than the regular instances may represent a departure from the WAN performance on other regular VMs as their fine-grained token bucket like mechanisms for resource sharing [16]. Recall that the observable bandwidth from a tenant’s perspective only occupies a small part of WAN capacity, so measuring variations in the underlying WAN can be more challenging; as such, there may be an opportunity to improve WAN bandwidth by detours when we consider data transfers on large-scale VMs.

Acknowledgements

Special thanks go to Zhe Wu and Muhammed Uluyol for collecting parts of our bandwidth and latency measurement datasets. This work was supported in part by National Science Foundation grants CNS-1563095, CNS-1463126, and CNS-1563849.

References

- [1] Raajay Viswanathan, Ganesh Ananthanarayanan, and Aditya Akella. Clarinet: Wan-aware optimization for analytics queries. In OSDI, 2016.
- [2] Q. Pu, G. Ananthanarayanan, P. Bodik, S. Kandula, A. Akella, V. Bahl, and I. Stoica. Low latency geo-distributed data analytics. In SIGCOMM, 2015.
- [3] K. Hsieh, A. Harlap, N. Vijaykumar, D. Konomis, G. Ganger, P. Gibbons, and O. Mutlu. Gaia: Geo-distributed machine learning approaching lan speeds. In NSDI, 2017.
- [4] V. Vulimiri, C. Curino, B. Godfrey, T. Jungblut, J. Padhye, and G. Varghese. Global analytics in the face of bandwidth and regulatory constraints. In NSDI, 2015.
- [5] Hajjat Mohammad, Liu Ruiqi, Chang Yiyang, Ng T. S. Eugene, and Rao Sanjay. Application-specific configuration selection in the cloud: impact of provider policy and potential of systematic testing. In INFOCOM, 2015.
- [6] Li Chen, Shuhao Liu, Baochun Li, and Bo Li. Scheduling jobs across geo-distributed datacenters with max-min fairness. In INFOCOM, 2017.
- [7] Hong Zhang, Kai Chen, Wei Bai, Dongsu Han, Tian Chen, Hao Wang, Haibing Guan, and Ming Zhang. Guaranteeing deadlines for inter-data center transfers. In EuroSys, 2016.
- [8] Zhiming Hu, Baochun Li, and Jun Luo. Flutter: Scheduling tasks closer to data across geo-distributed datacenters. In INFOCOM, 2016.
- [9] Shuhao Liu, Hao Wang, and Baochun Li. Optimizing shuffle in wide-area data analytics. In ICDCS, 2017.
- [10] David Andersen, Hari Balakrishnan, Frans Kaashoek, and Robert Morris. Resilient overlay networks. In SOSP, 2001.
- [11] Krishna P. Gummadi, Harsha V. Madhyastha, Steven D. Gribble, Henry M. Levy, and David Wetherall. Improving the reliability of internet paths with one-hop source routing. In OSDI, 2004.
- [12] Osama Haq, Mamoon Raja, and Fahad R. Dogar. Measuring and improving the reliability of wide-area cloud paths. In WWW, 2017.
- [13] Junchen J., Rajdeep D., Ganesh A., A.Philip C., N. Venkata P., Vyas S., Esbjorn D., Marcin G., Dalibor K., Renat V., and Hui Z. Via: Improving internet telephony call quality using predictive relay selection. In SIGCOMM, 2016.
- [14] Simon Peter, Umar Javed, Qiao Zhang, Doug Woos, Thomas Anderson, and Arvind Krishnamurthy. One tunnel is (often) enough. In SIGCOMM, 2014.
- [15] Omid Alipourfard, Hongqiang Harry Liu, Jianshu Chen, Shivaram Venkataraman, Minlan Yu, and Ming Zhang. Adaptively unearthing the best cloud configurations for big data analytics. In NSDI, 2017.
- [16] Neda Nasiriani, Cheng Wang, George Kesidis, and Bhuvan Uргаonkar. Using burstable instances in the public cloud: When and how? In SIGMETRICS, 2017.