

CloudPilot: Flow Acceleration in the Cloud

Kfir Toledo^{a,b}, David Breitgand^b, Dean Lorenz^b, Isaac Keslassy^a

^a*Technion, Haifa, Israel*

^b*IBM Research, Haifa, Israel*

Abstract

TCP-split proxies have been previously studied as an efficient mechanism to improve the rate of connections with large round trip times. These works focused on improving a single flow. In this paper, we investigate how strategically deploying TCP-split proxies in the cloud can improve the performance of geo-distributed applications entailing multiple flows interconnecting globally-distributed sources and destinations using different communication patterns, and being subject to budget limitations.

We present *CloudPilot*, a Kubernetes-based system that measures communication parameters across different cloud regions, and uses these measurements to deploy cloud proxies in optimized locations on multiple cloud providers. To this end, we model cloud proxy acceleration and define a novel *cloud-proxy placement problem*. Since this problem is NP-Hard, we suggest a few efficient heuristics to solve it. Finally, we find that our cloud-proxy optimization can improve flow completion time by an average of $3.6\times$ in four different use cases.

1. Introduction

Motivation. Over the last few years, the fierce competition among cloud providers has led them to spend billions on expanding their global presence by building data-centers worldwide and laying out high-speed lines to interconnect them [45, 6, 33]. Clients can now build on-demand cloud overlay networks comprising cloud nodes in different regions to route application traffic through the cloud rather than through the public internet [24, 2, 17, 34, 27, 3, 1, 31].

Several studies [34, 27, 3, 1, 31] show how we can increase the rate of a flow by using cloud-based TCP-split proxies. As Fig. 1 shows, this method splits a single TCP flow into several connections with shorter round trip time (RTT). By reducing the RTT for each connection, the overall transmission rate is improved.

While the above works focus on accelerating single flows, it is unclear how to strategically deploy a limited set of cloud-based TCP-split proxies to improve the performance of global geo-distributed applications, with sources and destinations that need to exchange large amounts of data. Such applications include distributed databases, batch file exchanges, VM migrations, and CDNs [42, 46, 35] (§3). The goal of this paper is to introduce CloudPilot, a Kubernetes-based system that is designed to optimize proxy placement and deploy the proxies to serve these applications.

Contributions. We make the following contributions.

- *Use cases.* We start by showing how geo-distributed applications can be modeled using four topological use-cases (§3).
- *CloudPilot.* We develop and deploy CloudPilot, a Kubernetes-based system that helps accelerate geo-distributed application traffic. Using communication parameters measured across different cloud regions, CloudPilot deploys TCP-split cloud proxies across multiple cloud providers to optimize the application transfer performance. We later present our proxy acceleration model for a single flow, and validate it using real-world CloudPilot-based measurement experiments (§4). The open-source CloudPilot code is available online [43].
- *Proxy placement optimization.* We explain how a natural metric of performance in geo-distributed applications is the *total Flow Completion Time (FCT)*, which we define to be the total time required to complete all necessary data transfers of the flow. We formally define a *cloud-proxy placement problem* to optimize this total FCT and prove that it is NP-hard (§5). Hence, we propose two families of heuristic algorithms to solve the problem: the *flow-greedy* algorithms, which greedily consider first the flows whose performance can most improve; and the *proxy-greedy* algorithms,

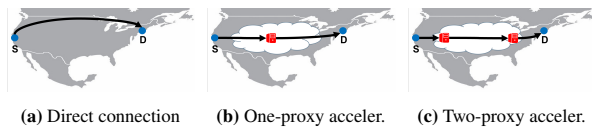


Figure 1: Connection types between host and destination: (a) direct connection, (b) one-proxy acceleration, and (c) two-proxy acceleration.

which greedily establish first the proxies that can most improve performance (§6).

- *Evaluations.* We evaluate our proposed heuristics both through extensive simulations with parameters measured from actual cloud providers, and through real-world CloudPilot-based cloud-proxy deployments. We find that our heuristic algorithms achieve significant FCT acceleration. For example, spending 50¢-per-flow to transfer 2GB-flows on Google cloud decreases the total FCT by factors of 2.7, 3.6, 3.9 and 4.3 for four different application use-cases. We also find that counter-intuitively, FCT acceleration significantly improves as last-mile bandwidth increases, especially beyond 100Mbps, heralding an increased impact for CloudPilot with the last-mile fiber-optics deployment (§7).

2. Related Work

TCP splitting. Many previous studies show the benefits of using TCP splitting over regular TCP [23, 30, 8, 21, 5, 38]. These works show how using TCP splitting proxies can improve the throughput in different environments such as mobile, web transfer over an HTTP connection, satellite connection with long distances, etc. We extend related work by focusing on placing the TCP proxies by demand in the cloud environment, exploiting the cloud infrastructure. Using the cloud environment, we can choose the proxy location, and our goal is to optimize the proxy location under cost limits.

Cloud infra-structure and performance. Companies like Amazon, Google, IBM, and Microsoft, spend significant effort and money into developing their clouds. They have a high-end infrastructure with optimized network algorithms. Each in-cloud provider can use highly optimized advanced protocols in its data centers. For example, Google uses TCP-BBR [4] and QUIC [9], and AWS uses SRD [40], which improves network traffic. The cloud providers build data centers worldwide that allow fast connections to users. In addition, cloud

providers make private internet connections [26] between data centers that will enable them to control the network routing efficiently and to pass only their own traffic. Exploiting these advantages, we can build a cloud overlay network with network acceleration (like TCP splitting) for geo-distributed applications.

Cloud overlay network. Several works show that forwarding using cloud proxies without TCP splitting capability provides little to no improvement [24, 1, 3]. Later research shows the benefit of using TCP splitting in the cloud overlay network. [3, 27, 34] show that using a single TCP splitting proxy can achieve up to 3× improvement over a direct internet connection. [3] also suggests using Multi Path-TCP (MP-TCP) to increase performance. However, MP-TCP is not always supported by communicating parties. [1, 31] show we can achieve better performance by using two TCP-splitting proxies with large buffers, one close to the source and one to the destination. In addition, they implement several improvements for the TCP splitting, like TCP turbo start, which can also be implemented in our system and further increase its performance. However, they only consider isolated flows and not a full system. Several recent works also analyze in more detail the performance of cloud communications [45, 6, 39, 33] and the additional benefits of cloud overlays [17].

Geo-distributed applications. Most research on the performance improvement of geo-distributed applications focuses on load-balancing mechanisms over direct TCP connections [47, 36, 22, 29, 19]. Our work complements these efforts by introducing TCP splitting, obtaining further significant performance gains.

Caching proxy placement. Several papers study the placement of cache proxies [16, 44] and HTTP-gathering proxies [2]. The TCP-split proxy placement problem is different. For example, in the above examples, it is preferable to place a cache as close to endpoints as possible, while in TCP splitting, the preferred location of a single proxy is in the middle.

3. Use-cases

We are interested in considering many data-intensive geo-distributed applications. Since they may vary considerably, to reason about them, we abstract away details and focus on their characteristic communication patterns, classifying them into four topological use-cases.

One-to-many. A single source broadcasts information to many destinations worldwide. Usually, a Content Delivery Network (CDN) will be used for most of the

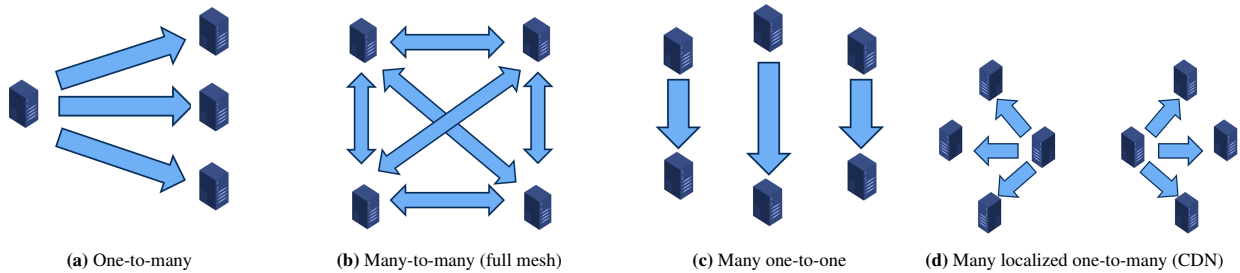


Figure 2: The topologies of four different use-cases.

destinations [35]. However, according to the Facebook statistics [10] some 1.8% of all live streams use direct connections due to cache misses. These remaining flows can be modeled using a star-like one-to-many pattern. Fig. 2(a) illustrates the one-to-many topology.

Many-to-many. Many nodes in different locations communicate in a full-mesh pattern, *e.g.*, in a geo-distributed database that transfers data between nodes to keep consistency [42]. Fig. 2(b) illustrates the many-to-many topology.

Many one-to-one. Topology with many unrelated source-destination pairs. One example is a VM migration application [46] that entails sending data from one data-center location to another for many unrelated VMs. Additional examples include backup between data centers, and file-sharing systems. Fig. 2(c) illustrates the many one-to-one topology.

Many localized one-to-many. Several sources that are geographically distributed and each broadcasts to many mostly-local destinations. One example is the traffic between CDN caches and their end-users [35]. Another is Twitch, an interactive live-streaming platform that offers three servers in three different continents [7]. Fig. 2(d) illustrates the many localized one-to-many topology.

In this paper, we focus on a static setting where applications have predictable traffic patterns (for instance, periodic backups), and we use this predictability to optimize the proxy placement by minimizing the total FCT of future flows. In future work, this could be extended to a speculative setting that uses various prediction models to estimate the upcoming future patterns. For example, an hourly process could consider the flow distribution in the last hour, then establish proxies for the next hour based on the expectation that the location distribution of future requests will be close enough. In addition, for simplicity, we focus in this paper on applications with at most dozens of flows per period.

4. Proxy Acceleration System

In this section, we introduce our CloudPilot system that utilizes TCP-split proxies to reduce FCT. We derive a model to predict the FCT based on proxy properties and network measurements and provide empirical evidence for the effectiveness of our approach. In the next section, we define the cloud-proxy placement problem that will be solved with our model.

4.1. CloudPilot

CloudPilot system. *CloudPilot* is a Kubernetes-based system that measures communication parameters across different cloud regions, and uses these measurements to deploy cloud proxies in optimized locations on multiple cloud providers. It is able to deploy new Kubernetes clusters on multiple cloud providers, deploy Kubernetes container instances, and connect between them. It also deploys iPerf3 measurement containers and HAProxy proxy containers. Finally, it implements the algorithms of this paper to decide where to deploy the proxies. The CloudPilot code is available online [43].

Deployment for FCT measurements. To obtain the FCT measurements below, CloudPilot spawns Kubernetes 1.22.2 clusters on demand for each pair of (source, destination) locations. The source cluster executes a Kubernetes container running an iPerf3 3.9 client [18] and the destination cluster executes a Kubernetes service with a backend container running an iPerf3 server. This way, FCT can be measured remotely between the source and destination over a direct cloud connection.

To obtain the FCT for connections with TCP-split proxies, CloudPilot creates additional Kubernetes clusters in different locations. CloudPilot can create three types of acceleration configurations: forwarding, one-proxy, and two-proxy.

For the forwarding acceleration, CloudPilot uses an Ubuntu 20.04 container. Traffic forwarding is done using appropriate iptables rules. For TCP splitting, CloudPilot uses an HAProxy 2.2.19 container that splits a

Table 1: Configuration of CloudPilot measurement experiments

| | source | dest. | 1 proxy | 2 proxies | |
|---|-----------------------------|---------------------------------|------------------------|---------------------------|----------------------|
| | | | | first | second |
| a | Israel <i>public net</i> | California <i>public net</i> | London <i>GCP</i> | London <i>GCP</i> | Oregon <i>GCP</i> |
| b | London <i>AWS</i> | Oregon <i>AWS</i> | Montreal <i>GCP</i> | London <i>GCP</i> | Oregon <i>GCP</i> |
| c | Ohio <i>AWS</i> | Mumbai <i>AWS</i> | London <i>GCP</i> | Virginia <i>GCP</i> | Mumbai <i>GCP</i> |
| d | S. Carolina <i>GCP</i> | Oregon <i>GCP</i> | Iowa <i>GCP</i> | S. Carolina <i>GCP</i> | Oregon <i>GCP</i> |

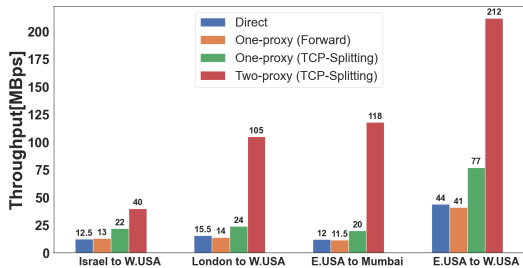


Figure 3: Rate comparison between (1) direct connection, (2) one-proxy forwarding, (3) one-proxy splitting and (4) two-proxy splitting, using different types of (source, destination, proxies) tuples as described in Table 1.

TCP connection into two connections with smaller RTT. In addition, when CloudPilot uses two-proxy splitting acceleration, it increases the TCP buffer size of the containers to match their intra-cloud bandwidths between the proxies.

FCT measurements. Fig. 3 presents the results of our preliminary measurement experiments intended to gain an intuition about different flow proxy acceleration options. The rationale is to understand which options are likely to result in most gains and focus our exploration on those settings. Table 1 summarizes the configuration of each experiment. Fig. 3 shows how a forwarding proxy obtains either negligibly better or even worse performance compared to a direct communication over public internet. This is consistent with previous studies [1, 27]. TCP-split proxies clearly outperform both direct connections and forwarding proxies. Therefore, in the remainder of this paper we focus on split proxies.

4.2. Proxy Acceleration Model

We now analyze how TCP-split proxies affect FCT and develop a model for estimating the FCT for several proxy deployment options: *direct connection*, *one proxy*, and *two proxies*, using measurable permanent properties of the end-hosts and the potential proxies. To do so, we consciously ignore the temporary impacts of the loss rate, queueing time, packet reorder-

ing, and similar effects along the packet path. In other words, we make the following simplifying assumptions: (1) each modeled flow transfers a large amount of data; (2) packet loss rate is negligible; (3) queueing time in the network is negligible vs. the propagation time; and (4) packet reordering is negligible. As we show in the next subsection, these assumptions are verified by our real cloud experiments.

Direct connection. We consider four main measurable factors affecting the FCT of a TCP flow from i to j .

Transfer size. Assume that i wants to transfer $\omega_{i,j}$ bits to j . Then FCT is directly proportional to the transfer size $\omega_{i,j}$ (using Assumption (1)).

Round Trip Time (RTT). The RTT equals $RTT_{i,j}$, its propagation component between i and j (Assumption (3)).

Maximum window size. Let $WND_{i,j}$ be the maximum possible window size between i and j , as limited by the respective OS configurations. Since we send at most $WND_{i,j}$ bytes per RTT, the rate between i and j is bounded by $\frac{WND_{i,j}}{RTT_{i,j}}$ [20].

Last-mile bandwidth. The flow’s rate is limited by both the last-mile egress bandwidth of source i and the last-mile ingress bandwidth of destination j . The last-mile bandwidth may reflect a variety of factors, including the internet service provider rate limit or the NIC speed. We denote as $BW_{i,j}$ the minimum of these two last-mile bandwidths.

Combining the above factors and applying assumptions, the flow rate is bounded by either $\frac{WND_{i,j}}{RTT_{i,j}}$ or $BW_{i,j}$, yielding an approximate rate of $R_{i,j} \approx \min\left(\frac{WND_{i,j}}{RTT_{i,j}}, BW_{i,j}\right)$. Its FCT $T_{i,j}^{\text{direct}}$ is approximated by $T_{i,j}^{\text{direct}} \approx \frac{\omega_{i,j}}{R_{i,j}}$.

One-proxy acceleration. FCT for a flow with one-proxy accelerator p is $T_{i,j}^p \approx \frac{\omega_{i,j}}{\min(R_{i,p}, R_{p,j})}$, because the flow rate is the rate of its slowest hop.

Two-proxy acceleration. FCT for a flow with two-proxy acceleration (p, q) is $T_{i,j}^{p,q} \approx \frac{\omega_{i,j}}{\min(R_{i,p}, R_{p,q}, R_{q,j})}$. As in [1], we assume that in the proxies, the maximum window sizes on the Internet side use the Linux default. However, on the internal cloud side they can be increased to take full advantage of the paid cloud bandwidth rates, namely $R_{p,q} \approx BW_{p,q}$, which is set by the cloud proxy capacity.

4.3. Model validation

Fig. 4 puts our model to test in the real world. It plots the real-world measured *rate* and *rate acceleration* against the predicted values using our model, for

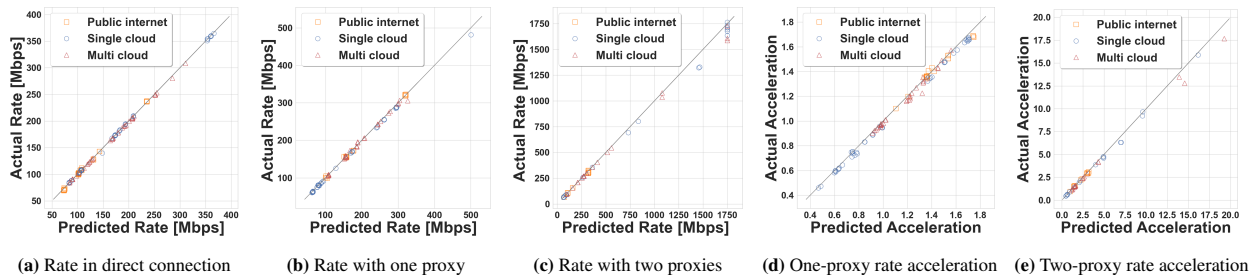


Figure 4: Model validation: Real-world *rate* vs. predicted rate using (a) direct connection; (b) one proxy; and (c) two proxies; then real-world *rate acceleration* (beyond direct connection) vs. predicted one using (d) one and (e) two proxies. In all cases, the model predictions seem close to real-world measured values.

all three types of connections. The figure shows sample results. Each data point averages 20 runs. Using CloudPilot we run three types of experiments: (1) *Public internet*, where source and destination are located in the public internet, and proxies are located in different regions of the same cloud, Google Cloud Platform (GCP) in this case. We use a desktop computer with Ubuntu v20.04 located at the Technion (Israel) as a host, and nine public iPerf3 servers around the world as destinations [18, 41]. (2) *Single cloud*, where the source, destination and proxies are sampled from 24 potential GCP locations. (3) *Multi cloud*, where source and destination are sampled from 9 IBM cloud locations, but the proxies are deployed in GCP. Overall we run over 75 different tuples (source, destination, proxies) with over 3,000 tests. The flow average rate is measured for 40 seconds. Our model prediction is based on a maximum window size of 2.875MB (observed default for Linux TCP) for all links, except in the cloud-facing links of the proxies where they are set to 500MB. The proxy bandwidth limitation is estimated as 1750Mb/s (the per-flow limitation of HAProxy), as it is tighter than the 2Gb/s link capacity of our used proxy machines.

5. Cloud-proxy placement problem

In this section, we present the cloud-proxy placement problem. First, we explain the problem informally to equip the reader with some intuition. Next, we introduce a formal notation and present a MILP (mixed-integer linear programming) formulation of the problem.

5.1. Informal Problem Definition

Given a set of source-destination pairs representing TCP flows, we want to find a feasible allocation of the flows to a set of TCP-split proxies in cloud regions, such that we minimize the total FCT (the sum of per-flow

FCTs). An allocation is feasible if (a) its cost is no greater than the overall predefined budget, and (b) for any proxy, the sum of bandwidth demands of all flows using this proxy is no greater than its capacity. Each flow can be allocated one, two, or zero proxies (the latter corresponds to a direct connection).

The cost of using a proxy comprises two components: (1) the proxy setup cost (*e.g.*, a virtual machine or a container with a specific bandwidth capacity), and (2) the data-transmission cost.

Note that (1) the setup cost is paid only once when multiple flows share a proxy, and (2) the cloud providers do not impose costs on the inbound network traffic. Only the outgoing traffic from the cloud proxy is billed (to different regions or to exit the cloud).

5.2. Problem Statement

Fig. 5 presents a formal MILP formulation for our cloud-proxy placement problem. Formally, we want to find a feasible assignment of proxies to flows such that the total FCT in the system is minimized, given constraints that reflect (1) the input sets and parameters presented on the top-left of Fig. 5, including the proxy setup costs, proxy data-transfer costs, and total budget; and (2) the per-flow FCT of each flow using any zero, one or two proxies, as computed by the proxy acceleration model of §4.

Theorem 1. *The cloud-proxy placement problem is NP-Hard.*

Proof. The 0/1 multiple-knapsack problem (MKP) [32] is a known NP-hard problem. In this problem, we need to place a subset of N non-splittable items in M bins. Each item i has a given positive weight w_i and profit p_i . The sum of the weights of all items in a bin j cannot exceed its capacity C_j . Our goal is to place a subset of items in the bins with maximum sum of the subset item

| Notation | Description |
|---|--|
| <i>Input Sets</i> | |
| \mathcal{S} | Set of all servers in the system, $s_i \in \mathcal{S}$ |
| \mathcal{F} | Set of all valid flows in the system, $f_{i,j} \in \mathcal{F}$ |
| \mathcal{L} | Set of all possible regions for proxies, $l \in \mathcal{L}$ |
| \mathcal{N} | Set of all possible instances for proxy, $n \in \mathcal{N}$ |
| \mathcal{P} | Set of all possible proxies $p \in \mathcal{P}$, \mathcal{P} contains all instances in all locations, $\mathcal{P} \equiv (\mathcal{L} \times \mathcal{N})$ |
| \mathcal{A} | Set of all possible proxy assignments (2 proxies, 1 proxy, or direct connection), $\mathcal{A} = (\mathcal{P} \times \mathcal{P}) \cup (\mathcal{P} \times \{0\}) \cup \{(0, 0)\}$ |
| <i>Input Parameters</i> | |
| $\omega_{i,j}$ | Data size to transfer by $f_{i,j}$ |
| $BW(p)$ | Bandwidth capacity of proxy $p \in \mathcal{P}$ |
| $C_{BW}(p)$ | Network traffic cost per Gigabyte using proxy $p \in \mathcal{P}$ |
| $C_{\text{setup}}(p)$ | Cloud proxy setup cost for proxy $p \in \mathcal{P}$ |
| B | Maximum allowed budget in the system |
| <i>Computed FCTs by CloudPilot (§4)</i> | |
| $T_{i,j}^{\text{direct}}$ | FCT of $f_{i,j}$ using direct path |
| $T_{i,j}^p$ | FCT of $f_{i,j}$ using one proxy $p \in \mathcal{P}$ |
| $T_{i,j}^{p,q}$ | FCT of $f_{i,j}$ using two proxies $p, q \in \mathcal{P}$ |
| <i>Decision Variables</i> | |
| $u_{i,j}^{p,q}$ | $\begin{cases} 1 & \text{flow } f_{i,j} \text{ uses proxies } p, q \in \mathcal{P} \\ 0 & \text{otherwise} \end{cases}$ |
| x_k | $\begin{cases} 1 & \text{if proxy } k \in \mathcal{P} \text{ is used} \\ 0 & \text{otherwise} \end{cases}$ |

| Input | |
|-------------------------|---|
| $\mathbf{T}_{ij}[p, q]$ | $= \begin{cases} T_{i,j}^{\text{direct}} & \text{if } (p, q) = (0, 0) \\ T_{i,j}^p & \text{if } (p, q) = (p, 0) \\ T_{i,j}^{p,q} & \text{otherwise.} \end{cases} \quad \forall (p,q) \in \mathcal{A}$ (FCTs from §4) |
| Optimization goal | |
| minimize | $\sum_{\forall f_{i,j} \in \mathcal{F}} \sum_{(p,q) \in \mathcal{A}} \mathbf{T}_{i,j}[p, q] \cdot u_{i,j}^{p,q} \quad (\text{total FCT})$ |
| Constraints | |
| | $\sum_{(p,q) \in \mathcal{A}} u_{i,j}^{p,q} = 1 \quad \forall f_{i,j} \in \mathcal{F} \quad (\text{one allocation per flow})$ |
| | $BW(k) \geq \sum_{\forall f_{i,j} \in \mathcal{F}} \omega_{i,j} \left(\sum_{(p,k) \in \mathcal{A}} \frac{u_{i,j}^{p,k}}{\mathbf{T}_{i,j}[p, k]} + \sum_{(k,q) \in \mathcal{A}} \frac{u_{i,j}^{k,q}}{\mathbf{T}_{i,j}[k, q]} \right) \quad \forall k \in \mathcal{P}$ (proxy capacity fulfills bandwidth demand) |
| | $x_k \leq \sum_{f_{i,j} \in \mathcal{F}, (p,q) \in \mathcal{A} \text{ s.t. } p=k \vee q=k} u_{i,j}^{p,q} \quad \forall k \in \mathcal{P}$ (0 if k unneeded) |
| | $ \mathcal{N} x_k \geq \sum_{f_{i,j} \in \mathcal{F}, (p,q) \in \mathcal{A} \text{ s.t. } p=k \vee q=k} u_{i,j}^{p,q} \quad \forall k \in \mathcal{P}$ (1 if k needed) |
| | $x_{k_1} = x_{(l, n_1)} \geq x_{(l, n_2)} = x_{k_2} \quad \forall n_1 < n_2 \in \mathcal{N}, \forall l \in \mathcal{L}$ (for proxies k_1, k_2 with the same location l , prefer the smaller index) |
| | $C_{\text{setup}}^{\text{total}} = \sum_{k \in \mathcal{P}} C_{\text{setup}}(k) x_k \quad (\text{total setup cost})$ |
| | $C_{\text{BW}}^{\text{total}} = \sum_{f_{i,j} \in \mathcal{F}} \omega_{i,j} \sum_{p,q \in \mathcal{A}} u_{i,j}^{p,q} (C_{\text{BW}}(p) + C_{\text{BW}}(q)) \quad (\text{total BW cost})$ |
| | $B \geq C_{\text{setup}}^{\text{total}} + C_{\text{BW}}^{\text{total}} \quad (\text{budget limitation})$ |

Figure 5: MILP formulation of cloud-proxy placement problem, with a table of used notations on the left.



Figure 6: Intuition for algorithm choices. (a) Flow-greedy algorithms pick the best proxy acceleration for the flow that benefits most, even if expensive, and tend to prefer two-proxy acceleration; while (b) proxy-greedy algorithms choose the single proxy that can most benefit the system by serving several flows, thus spending the budget more efficiently.

profits. Given any 0/1 MKP instance, we define a corresponding instance of the cloud-proxy placement problem with a single proxy location, and show that solving it would also solve the 0/1 MKP problem. We define N flows, and can freely choose their flow rates as w_i (we can arbitrarily change the maximum window, given an infinite BW and a fixed RTT), and flow FCT gain (difference between FCT in direct connection and FCT using the proxy) as p_i (we can arbitrarily change the data size of flow i). We define the budget as B . We set the bandwidth cost as $C_{BW} = 0$ and proxy-setup cost as $C_{setup} = \frac{B}{M}$, so the budget allows exactly M proxy instances at this location. We set the bandwidth capacity of proxy j to C_j . Since there is only one proxy location, using two-proxy acceleration is never beneficial.

If there is a solution to our problem, we can also solve the 0/1 MKP. Hence, by reducing the 0/1 MKP to the above problem, we find it is NP-Hard. \square

6. Algorithms

Since the cloud-proxy placement problem is NP-hard, we propose two families of greedy approximation algorithms:

- The *flow-greedy* family of algorithms, where we greedily allocate flows, one at a time.
- The *proxy-greedy* family of algorithms, where we greedily allocate proxies, one at a time.

Fig. 6 provides an example for understanding the intuition behind the two families of algorithms.

6.1. Flow-greedy algorithms

We propose two versions of the flow-greedy algorithm that differ only by the gain calculation, namely the order of processing flows.

F-FCT (Flow-greedy FCT). The pseudo-code for *F-FCT* is given in Alg. 1. It takes as input (Line 1) the set of flows, the set of proxies, and the overall budget. For

Algorithm 1 Flow Greedy FCT (F-FCT)

```

1: MAIN( $\mathcal{F}^0, \mathcal{P}^0, B^0$ )  $\triangleright$  Flow and proxy sets, budget
2:  $\mathcal{D} \leftarrow (\mathcal{L} \times \mathcal{L}) \cup (\mathcal{L} \times \{0\}) \cup (0, 0)$   $\triangleright \mathcal{D}$  is a list of every possible proxy allocation location
3: for  $f \in \mathcal{F}^0$  do
4:    $a_f \leftarrow (0, 0)$   $\triangleright a_f$  is the allocation for  $f$ , initialized as a direct connection
5:    $\mathbf{G}_f \leftarrow \mathcal{D}$ , sorted non-increasing by  $\text{GAIN}(f, d) \forall d \in \mathcal{D}$   $\triangleright \mathbf{G}_f$  is a list of all possible  $\mathcal{A}$  for  $f$  sorted by gain
6:    $r_p \leftarrow BW(p) \quad \forall p \in \mathcal{P}^0$   $\triangleright$  Available proxy bandwidth
7:    $\mathcal{P} \leftarrow \emptyset$   $\triangleright$  Allocated proxies so far
8:    $B \leftarrow B^0$   $\triangleright$  Remaining budget
— end of initialization —
9:    $\mathcal{F} \leftarrow \mathcal{F}^0$   $\triangleright$  Flows without allocated proxies
10:  while  $\mathcal{F} \neq \emptyset$  do GREEDY-STEP
11:  return  $\sum_{f \in \mathcal{F}^0} \mathbf{T}_f[a_f], \{a_f\}_{f \in \mathcal{F}^0}$   $\triangleright$  Return overall score and allocation per flow
12: GAIN( $f, a$ )  $\triangleright$  FCT reduction for  $f$  with allocation  $a = (p, q)$ 
13: return  $\mathbf{T}_f[0, 0] - \mathbf{T}_f[a]$ 
14: GREEDY-STEP
15:    $f \leftarrow \arg \max_{f \in \mathcal{F}} (\text{GAIN}(f, \mathbf{G}_f.head))$   $\triangleright$  Greedily choose flow
16:    $a, b \leftarrow \text{FIND-PROXY-INSTANCES}(f)$ 
17:   if  $b \leq B$  then ALLOCATE( $f, a, b$ )
18: FIND-PROXY-INSTANCES( $f$ )
19:    $(l_1, l_2) \leftarrow \mathbf{G}_f.head$ 
20:    $r \leftarrow \omega_f / \mathbf{T}_f[l_1, l_2]$   $\triangleright$  Flow BW requirement
21:    $a \leftarrow (\text{FIND-PROXY-AT}(l_1, r), \text{FIND-PROXY-AT}(l_2, r))$ 
22:    $b \leftarrow \sum_{p \in a, p \neq 0} \omega_f C_{BW}(p) + \sum_{p \in a, p \notin \mathcal{P}} C_{setup}(p)$   $\triangleright$  Marginal allocation cost
23:   return  $a, b$   $\triangleright$  Return chosen proxy instances,  $a$ , and their marginal cost,  $b$ 
24: FIND-PROXY-AT( $l, r$ )  $\triangleright$  Find proxy at location  $l$  with  $r$  free capacity
25:    $\mathcal{P}_l \leftarrow \{p \in \mathcal{P} \text{ s.t. } p\text{'s location is } l\}$   $\triangleright \mathcal{P}_l \leftarrow \emptyset$  if  $l = 0$ 
26:   if  $\exists p \in \mathcal{P}_l$  s.t.  $r_p \geq r$  then
27:     return  $p$   $\triangleright$  Proxy instance  $p$  has enough capacity for the flow
28:   if  $\exists p \in \mathcal{P}^0 \setminus \mathcal{P}$  s.t.  $p$ 's location is  $l$  then  $\triangleright$  Always False if  $l = 0$ 
29:     return  $p$   $\triangleright$  New proxy instance
30:   return 0
31: ALLOCATE( $f, a, b$ )  $\triangleright$  Allocate  $a$  to  $f$  with budget  $b$ 
32:    $a_f \leftarrow a, B \leftarrow B - b, \mathcal{F} \leftarrow \mathcal{F} \setminus \{f\}$ 
33:   for  $p \in a$  do
34:      $\mathcal{P} \leftarrow \mathcal{P} \cup \{p\}$   $\triangleright$  Add proxy if new
35:      $r_p \leftarrow r_p - \omega_f / \mathbf{T}_f[a]$   $\triangleright$  Update available capacity

```

each flow, *F-FCT* initializes the allocation to a direct connection (Line 4) and computes its gain for every possible proxy allocation (Line 5). That is, it considers all possible locations for one proxy or one proxy pair and computes the gain w.r.t. a direct connection. The possible allocations for each flow are sorted by their gain.

The main loop (Line 10) greedily processes flows one at a time, launching the greedy function that examines

the best proxy locations for each flow, and updates the flow with the highest gain (Line 15). It then finds concrete proxy instances at these locations and calculates the cost of allocating these instances to the flow. If the budget allows, then the allocation for the flow is completed (Line 31) and the greedy step concludes. In order to find concrete proxy instances and their marginal cost (Line 18), Alg. *F-FCT* calculates the needed rate through the allocated proxies (Line 20). Then, it looks for a proxy with enough free capacity at each location (Line 24). It first tries to find an existing proxy with enough available capacity (Lines 25-26); if that fails, it uses a new proxy instance. The cost of the allocation (Line 22) includes the bandwidth cost of each proxy and the setup cost if a new proxy instance is required.

The gain function used by Alg. *F-FCT* (Line 13) only considers the reduction in FCT, regardless of its impact on the total budget. As illustrated in Fig. 6, *F-FCT* prefers expensive two-proxy acceleration types that strongly reduce the FCT, rather than cost-efficient one-proxy connections. Intuitively, *F-FCT* is best to use when the budget is nearly unlimited.

Time complexity. Let m be the number of flows, L the number of regions, and n the number of instances for each region. The size of \mathcal{D} is $O(L^2)$, so sorting for all flows requires $O(mL^2 \log(L))$ time. Each greedy step requires $O(m+n)$ time, $O(m)$ to find the best flow¹ and $O(n)$ to find its concrete proxy allocation. One flow is removed after each successful greedy step, thus the total time for the successful steps is $O(m(m+n)) = O(m^2)$ (since $n \leq m$). In order to bound the work required to process unsuccessful greedy steps, after each successful allocation, we make sure that the best potential allocation for each flow (the head of its sorted list) falls within the remaining budget. This is done by removing infeasible allocations from the head of each flow's list (in $O(1)$ per removal).² There are at most $O(mL^2)$ such removals, so overall $O(mL^2)$ time is required. Summing all, we get $O(m(L^2 \log(L) + m))$ time for Alg. *F-FCT*.

F-Cost (Flow-greedy FCT per cost). This algorithm is similar to the previous one, but considers cost when greedily choosing flows to process. The pseudo code is given in Alg. 2; it uses the same code of Alg. 1 with the GAIN function replaced. The idea is to scale down the gain for each allocation by its expected cost. The flow rate for each potential allocation can be computed at initialization from its expected FCT, so the BW cost of

each allocation is known. However, the exact setup cost for each allocation cannot be known at initialization, since it depends on whether the allocation would use an existing proxy with enough free capacity or would require a new instance. Instead, Alg. 2 attributes a fraction of the setup cost for every allocated flow using the ratio of the flow rate to the capacity of the proxy. Note that this cost-based gain is only used to sort the flows and allocations and is *not* used to calculate the actual allocation cost (Alg. 1, Line 22). With this algorithm, we get better performance under a limited budget. The time complexity is the same as for Alg. 1.

Algorithm 2 Flow Greedy Cost (F-Cost) *extends* Alg. 1

```

1: GAIN(f, a)
2:   return  $\frac{\mathbf{T}_f[0,0] - \mathbf{T}_f[a]}{\sum_{p \in a} (\omega_f C_{BW}(p) + C_{\text{setup}}(p) \frac{\omega_f / \mathbf{T}_f[a]}{BW(p)})}$   $\triangleright \frac{\omega_f}{\mathbf{T}_f[a]}$  is
   f's rate

```

6.2. Proxy-greedy algorithms

In the *proxy-greedy* family of algorithms, we choose the best proxy locations incrementally in a greedy manner. We start from an empty proxy set and add a few proxies at a time, so long as the overall FCT improves. At each greedy step, we generate a list of candidate proxy sets, and choose the one with the best total FCT. The total FCT for each candidate proxy set is computed using Alg. 1. The difference between the algorithms is in the way the candidate sets are generated at each greedy step.

1-P (one-proxy greedy). This is the basic proxy-greedy algorithm, its pseudo-code is given in Alg. 3. Alg. 3 creates a candidate set that includes all possibilities of adding a single proxy instance to the existing set (Lines 5&14). For every possible location in \mathcal{L} , it creates a candidate proxy set that includes the existing proxies plus a new proxy instance at that location. Then, at each greedy step, Alg. 1 is called for every proxy set in the candidate set to compute its FCT score (Line 6).³ If a candidate has a better total FCT score then its allocation is saved. The algorithm returns if no candidate proxy set improves the total FCT (Line 10).

Time complexity The candidate set size is bounded by the number of locations and the number of greedy steps is $O(m)$, since there are at most 2 proxies per flow. Thus there are (mL) calls to Alg. 1. Each call requires $O(m(L^2 \log(L) + m))$, however the initialization sorting can be cached to reduce subsequent calls to

¹Note that the gain for each allocation can be cached.

²This implementation detail is omitted from Alg. 1 to simplify the presentation. The check can be done in $O(1)$ by caching $\max_{p \in \mathcal{P}_i} r_p$ (Line 26).

³Note that the proxy set defines how many instances are available at each location, thus FIND-PROXY-AT may return 0 also for $l \neq 0$.

Algorithm 3 One-Proxy Greedy (1-P)

$\mathcal{Q}^{\mathcal{F}} = \{a_f\}_{f \in \mathcal{F}}$ denotes a set of proxy allocations a_f for all flows f

```
1: MAIN
2:  $Score_{min}, \mathcal{Q}_{min}^{\mathcal{F}}, \mathcal{P}_{min} \leftarrow \infty, \{(0, 0)\}_{f \in \mathcal{F}}, \emptyset$   $\triangleright$  Init
3: do
4:    $update \leftarrow False$   $\triangleright$  Flag indicates score improvement
5:   for each  $\mathcal{P}$  in CANDIDATE-PROXY-SETS( $\mathcal{P}_{min}$ ) do
6:      $Score, \mathcal{Q}^{\mathcal{F}} \leftarrow$  FLOW-GREEDY-FCT( $\mathcal{F}, \mathcal{P}_l, B$ )
7:     if  $Score < Score_{min}$  then
8:        $Score_{min}, \mathcal{Q}_{min}^{\mathcal{F}}, \mathcal{P}_{min} \leftarrow Score, \mathcal{Q}^{\mathcal{F}}, \mathcal{P}_l$ 
9:        $update \leftarrow True$ 
10:  while  $update$   $\triangleright$  Stop if no candidate improved score
11:  return  $Score_{min}, \mathcal{Q}_{min}^{\mathcal{F}}$ 
```

```
12: CANDIDATE-PROXY-SETS( $\mathcal{P}$ )
13:  return ADD-ONE-PROXY( $\mathcal{P}$ )
```

```
14: ADD-ONE-PROXY( $\mathcal{P}$ )
15:  for each region  $l \in \mathcal{L}$  do
16:    choose a proxy  $p_l$  in region  $l$  s.t.  $p_l \notin \mathcal{P}$ 
17:     $\mathcal{P}_l \leftarrow \mathcal{P} \cup \{p_l\}$ 
18:  return  $\{\mathcal{P}_l\}_{l \in \mathcal{L}}$ 
```

Algorithm 4 Two-Proxy Greedy (2-P) *extends* Alg. 3

```
1: CANDIDATE-PROXY-SETS( $\mathcal{P}$ )
2:  return ADD-TWO-PROXIES( $\mathcal{P}$ )
```

```
3: ADD-TWO-PROXIES( $\mathcal{P}$ )
4:   $\mathbb{P}^2 \leftarrow \emptyset$ 
5:   $\mathbb{P}^1 \leftarrow$  ADD-ONE-PROXY( $\mathcal{P}$ )
6:  for each combination  $\mathcal{P}^l \in \mathbb{P}^1$  do
7:    append ADD-ONE-PROXY( $\mathcal{P}^l$ ) to  $\mathbb{P}^2$ 
8:  return  $\mathbb{P}^2$ 
```

$O(m(L^2 + m))$.⁴ The overall time complexity is thus $O(m^2L(L^2 + m))$. Let P denote the number of proxies returned by the algorithm. Both the number of greedy steps and the number of available locations for each is bounded by P . Thus there are (PL) calls to Alg. 1 each requiring $O(m(P^2 + m))$. The overall complexity becomes $O(mPL(P^2 + m))$, which is tighter in practice as P is limited by the overall budget.

2-P (two-proxy greedy). The algorithm is based on Alg. 3, but with a candidate set that now includes all possibilities of adding two-proxy instances to the existing set (Line 3). The implementation reuses ADD-ONE-PROXY to generate the candidate set. The candidate set size is now $O(L^2)$, therefore the overall time complexity increases to $O(m^2L^2(L^2 + m))$ and $O(mPL^2(P^2 + m))$.

2-P RB (two-proxy greedy with rollback). The algorithm is again based on the Alg. 3, but with a candidate

⁴In practice, the bound on L is smaller for most calls, as we only need to consider the locations that are covered by each particular candidate proxy set $\{L \text{ s.t. } \mathcal{P}_l \neq \emptyset\}$.

Algorithm 5 Two-Proxy Rollback (2-P RB) *extends* Alg. 3

```
1: CANDIDATE-PROXY-SETS( $\mathcal{P}$ )
2:   $\mathbb{P}^{RB} \leftarrow \emptyset$ 
3:  for each  $p \in \mathcal{P}$  do
4:    append ADD-TWO-PROXIES( $\mathcal{P} \setminus p$ ) to  $\mathbb{P}^{RB}$ 
5:  return  $\mathbb{P}^{RB}$ 
```

set that now includes all possibilities of *removing* one proxy and adding two proxy instances to the existing set (Line 1). The idea is to avoid local minima by allowing the greedy algorithm to rollback one of the existing proxy allocations when it adds new proxies. Here we reuse ADD-TWO-PROXY from Alg. 4. Now the size of the candidate set is $O(L^3)$, so the overall time complexity is $O(m^2L^3(L^2 + m))$ or $O(mPL^3(P^2 + m))$. Although the above theoretical complexity bound is high, we found the actual run-time to be acceptable in practice. Both the number of world-wide cloud geographic locations and the number of flows is relatively small (dozens). Due to budget constraints, the number of allocated proxies is even smaller.

7. Evaluation

First, in simulations based on real-world parameters, we study the impact of several key model parameters and evaluate the performance of our algorithms on the use cases of §3. Then, in CloudPilot-based real-world cloud-environment experiments using Kubernetes and HAProxy, we confirm that the model predictions are close to reality, and that the proxy acceleration can be significant.

7.1. Settings

Runs. Each simulation data point is an average of 30 runs.

Proxy locations. We use 18 actual GCP regions for possible locations of the cloud proxies. The RTTs between the proxies are measured by CloudPilot and are consistent with a GCP RTT benchmark [11]. Due to lack of space, we present only the GCP results, but we obtained similar results in other cloud platforms that we checked, *e.g.*, IBM cloud.

Source and destination locations. To deploy each source, we first choose a random proxy, then select a location such that it has a reasonably small RTT to this proxy. We randomly choose locations with $RTT_{i,p} < 40\text{ms}$, corresponding to some 8,000Km using optical fibers [28]. Destination locations are chosen in the same way.

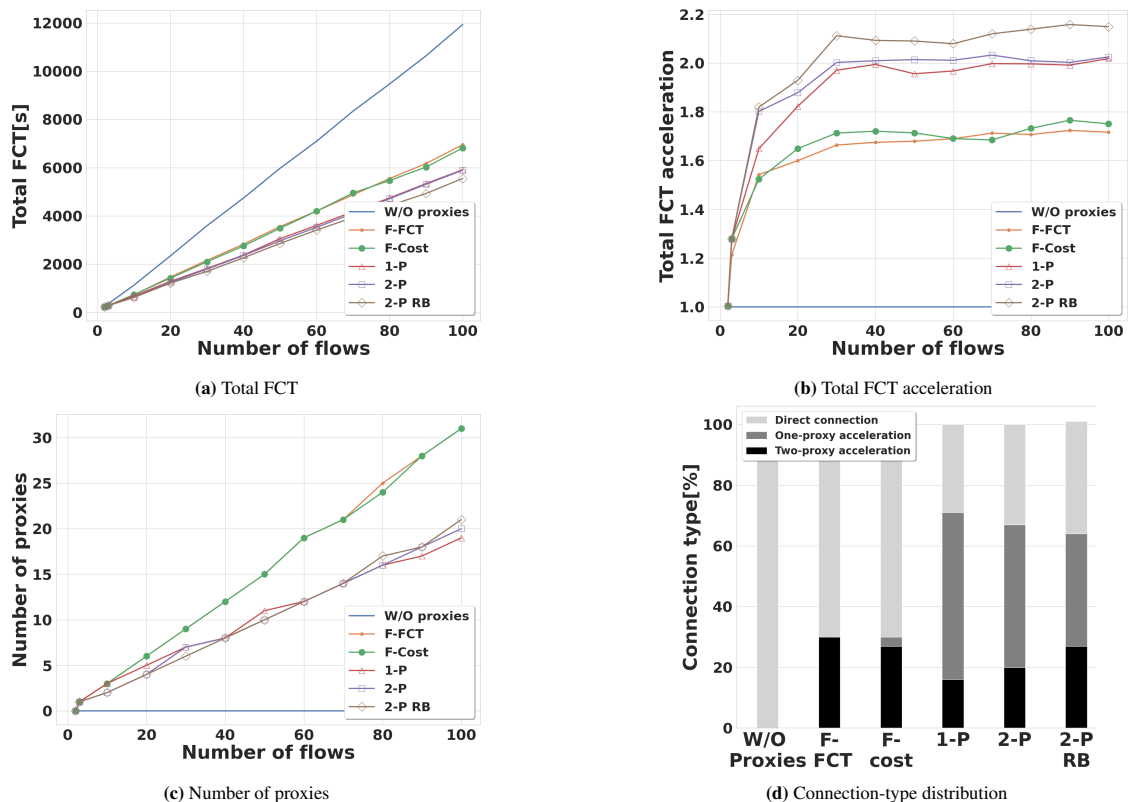


Figure 7: Impact of number of flows in one-to-many use case. (a) shows that the sum of flow FCTs tends to grow linearly with the number of flows in all algorithms, *i.e.*, the average FCT tends to converge, albeit the flow-greedy family achieves worse results than the proxy-greedy family. (b) illustrates the FCT acceleration when compared to a baseline without proxy. The proxy-greedy family of algorithms outperforms the flow-greedy algorithms and obtains over $2\times$ acceleration. (c) shows that the flow-greedy algorithms tend to spend a larger share of the budget on establishing proxies, while proxy-greedy algorithms focus on sharing proxies and spending on bandwidth. (d) details each family’s connection type distribution with 60 flows, confirming the intuition from Fig. 6 that flow-greedy algorithms tend to choose expensive single-use two-proxy accelerations for flows, while proxy-greedy algorithms prefer cheaper one-proxy accelerations with proxy sharing.

Network parameters. We set the transferred data size as $\omega = 2\text{GB}$ for each flow. We use a default constant proxy setup cost of $C_{\text{setup}}(p) = 50\text{¢}$ and constant bandwidth cost of $C_{\text{BW}}(p) = 8\text{¢}$ per GB for all proxies and regions, approximating the GCP prices [15, 14]. We set the proxy bandwidth capacity to 2Gbits since this is a standard egress bandwidth of a container on GCP [13]. We set the last-mile bandwidth BW of all our end-hosts to be 1Gbps, planning for a next-generation widespread gigabit access, at least among corporate customers [25]; except for the multi-flow servers, such as in the CDN and one-to-many use cases, which are not constrained by last-mile bandwidth. We use the default Linux window size for all servers and for one-proxy connections. For two-proxy connections, we increase the window size to 500MB for intra-cloud communications only.

7.2. System Evaluations

Impact of number of flows. We start by evaluating the impact of the number of flows on performance in a one-to-many use case. One source in Tokyo transfers data to each destination. At each step, we increment the number of flows by randomly adding a new destination worldwide. We set the budget proportionally to the number of flows. Fig. 7(a) and Fig. 7(b) show that proxy-greedy algorithms improve the total FCT in the system and outperform the flow-greedy algorithms. Fig. 7(c) illustrates how proxy-greedy algorithms use less proxies. Then, Fig. 7(d) shows that this is because proxy-greedy algorithms prefer having many flows share a single proxy for one-proxy acceleration. By saving on the proxy setup cost, they can accelerate more additional flows. By contrast, flow-greedy algorithms rely on expensive two-proxy acceleration.

Impact of budget. Fig. 8 shows the influence of bud-

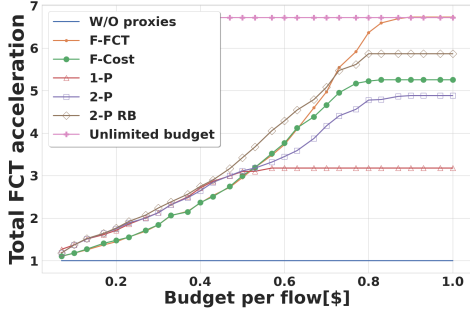


Figure 8: FCT acceleration for a one-to-many topology with 30 flows and different budgets. For instance, for a total budget of 14.5\$, *i.e.*, 50¢ per flow, it is possible to achieve some $3.1\times$ improvement in the average FCT. The P incr family better leverages low budgets and achieves higher accelerations. On the other hand, with high budgets, the F-greedy FCT that tends to pick the best and most expensive two-proxy acceleration choices manages to achieve the unlimited-budget bound, while the other algorithms cannot improve their choices that were picked greedily.

get on the overall FCT improvement for each algorithm. The FCT improvement is compared to the system performance without any proxy acceleration. In this simulation, we keep the settings of the previous one-to-many evaluation and consider 30 flows. We increase the total budget by 1\$ at every stage of the test. We can see that the proxy-greedy algorithms are superior for low and medium budgets. The gap between the F-FCT algorithm and the proxy-greedy algorithms is reduced for larger budgets. In addition, with a high budget, the F-FCT algorithm gets the best result, because it always picks the best proxy locations regardless of cost. Eventually, it gets the same improvement as in the case of an unlimited budget (when we use the best proxy locations for each flow without budget concerns). Consequently, when there is no budget limit, we should look at each flow separately and see that the two-proxy acceleration is the best option.

Impact of cost parameters. Fig. 9 and 10 show the impact of cost parameters. Each boxplot box represents the results between the 25th and 75th percentiles of 30 runs.

Data-transfer cost. Fig. 9 doubles the data-transfer cost to 16¢ per GB and zeroes the proxy-setup cost. In this case, almost all the algorithms reach the same results (Fig. 9(a)). The two-proxy acceleration method is preferred (Fig. 9(b)) because the cost of an extra proxy is neglected. So, we can create different proxies for each flow without utilizing and combining flows to the same proxies. This example applies when we transfer vast amounts of data between several flows and the data transfer price is dominant. In addition, the 1-P algo-

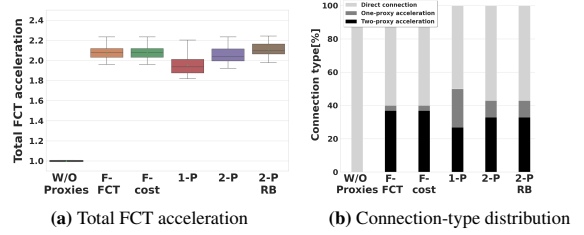


Figure 9: Impact of data-transfer cost. (a) When the data-transfer costs are dominant, both families get similar acceleration results. (b) When there is no cost for setting proxies, all algorithms use two-proxy acceleration.

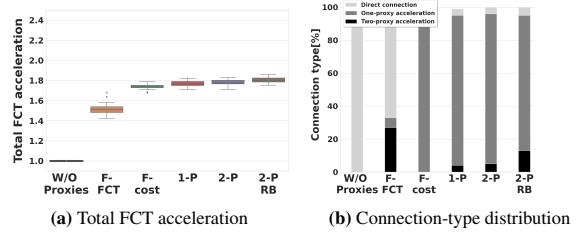


Figure 10: Impact of proxy setup cost. (a) With dominant proxy-setup cost, the cost-efficient algorithms obtain better accelerations. (b) In this case, the algorithms use one-proxy acceleration and reduce the number of proxies in the system.

rithm receives the lowest results because it is less capable of efficiently placing the corresponding two proxies.

Proxy-setup cost. Fig. 10 zeroes the data-transfer cost and doubles the proxy-setup cost to 1\$. Cost-efficient algorithms that share proxies obtain better results (Fig. 10(a)). Because we want to limit as much as possible the number of proxies, in this case, the F-greedy cost and the proxy-greedy families get the best results by choosing more cost-efficient flows and display a larger usage of the one-proxy acceleration (Fig. 10(b)). This example applies when every proxy setup has a large overhead (cost) in the organization.

Impact of file data size. Fig. 11 shows the impact of the file data size on the total acceleration. We use the F-FCT greedy algorithm, assuming 50¢ per flow. As we can see, we cannot accelerate files of small data size since the bandwidth of the source is sufficient to send the small files and there is no need to use TCP-split proxies. The FCT acceleration increases as we increase the file data size, and more flows use the TCP-split proxies. In addition, the data-transfer cost increases when we increase the file data size. We can accelerate fewer flows with a large data size for a fixed budget, therefore the acceleration is decreasing. On the other hand, for an incremental budget per file data size, the budget and ac-

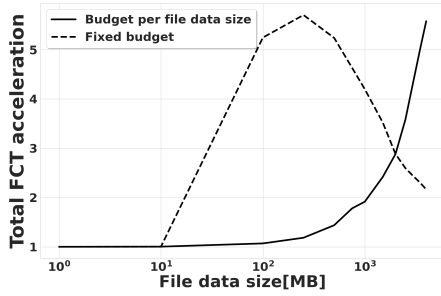


Figure 11: Impact of file data size: We do not get FCT improvements for small files. As we increase the file data size, the FCT acceleration improves. For a fixed budget, when the data transfer cost becomes dominant, the acceleration decreases.

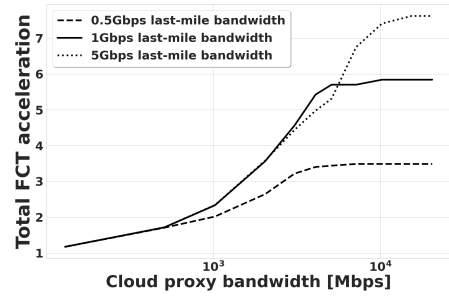


Figure 13: Impact of cloud proxy bandwidth: As the cloud proxy bandwidth grows, FCT acceleration increases, especially when the cloud-proxy capacity exceeds the last-mile bandwidth.

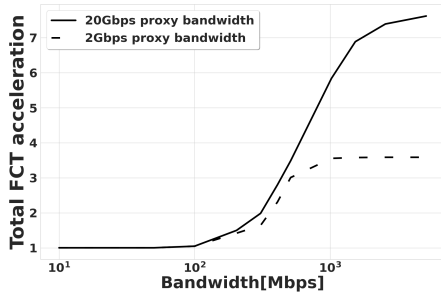


Figure 12: Impact of last-mile access bandwidth: As it grows, FCT acceleration increases significantly, especially when the cloud-proxy capacity is 20Gbps.

celeration increase when we increase the file data size.

Impact of last-mile access bandwidth. Fig. 12 shows the impact of the last-mile access bandwidth on the total acceleration. We use the 2-P RB algorithm, assume 50¢ per flow, and compare two types of proxy: small (2Gbps), as in current cheapest proxies, and large (20Gbps), assuming next-generation proxies will have larger limits. As the access bandwidth grows beyond some 100Mbps, the FCT acceleration increases significantly, especially with the large proxy capacity. This is because the flows are less constrained by the last-mile bandwidth, but rather by the long RTT, in which case cloud proxies with TCP splitting help more. This may partly justify the current increased interest in cloud proxies, as last-mile fiber-optics deployment becomes wider.

Impact of cloud proxy bandwidth. Fig. 13 shows the impact of the proxy cloud bandwidth on the total acceleration. We use the 2-P RB algorithm, assume 50¢ per flow, and compare three types of last-mile bandwidth: small (500Mbps), medium (1Gbps), and large (5Gbps). As cloud proxy bandwidth grows, accelera-

tion increases for all cases. The main reason for the acceleration increase is that with cloud proxies using larger bandwidth, we can aggregate more flows in the same proxy and reduce the proxy set-up cost in the system. The acceleration grows until we reach the budget limit. In addition, we can see that with a larger last-mile bandwidth, we can better utilize the larger-bandwidth cloud proxies. Interestingly, in some cases, the 1Gbps last mile-bandwidth FCT acceleration is greater than the 5Gbps last-mile bandwidth acceleration. These cases can happen only where the cloud proxy bandwidth is smaller than the last-mile bandwidth (5Gbps). In those cases, for the smaller last-mile bandwidth, the total FCT does not improve, but the FCT acceleration does. So, we can get acceleration even if the cloud proxy’s bandwidth is smaller than the last-mile bandwidth, because we reduce the RTT of the connection with the TCP split. Still, to get higher accelerations, it is better to use a cloud proxy with larger bandwidth than the sources and destinations.

Use cases. Fig. 14 shows the algorithm performance results for all four different use cases of §3. In all cases, we assume a 50¢ budget per flow and measure the total-FCT improvement for 60 flows. In the *many localized one-to-many* CDN-like topology, we first randomly select three sources in three different continents: Asia, Europe and North America. Since CDNs are not perfect, at each step, when we sample a random destination, it connects to its closest source with 90% probability, to its second-closest source with 7% probability, and to its farthest source with 3%. The first three use cases get high accelerations. In the fourth, the acceleration is smaller due to the shorter average distance, but non-negligible due to the many available proxy locations that enable us to perform a two-proxy acceleration (especially with F-FCT).

Run-time. The simulation time for an example run

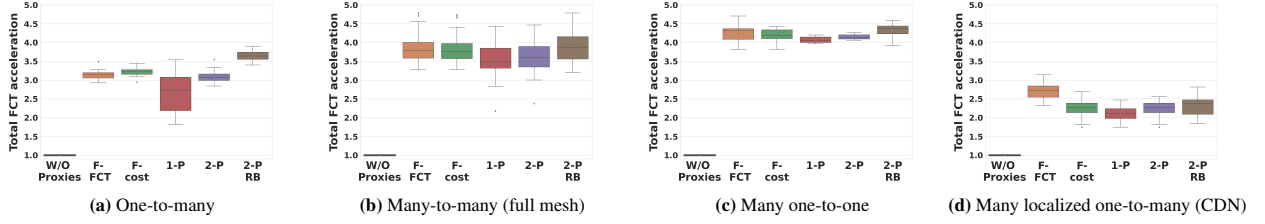


Figure 14: Impact of use cases with 50¢ per-flow budget. As expected, we get a strong acceleration for the first three use cases. In the fourth case that exploits CDN localization, smaller distances enable less proxy options and therefore a lower acceleration.

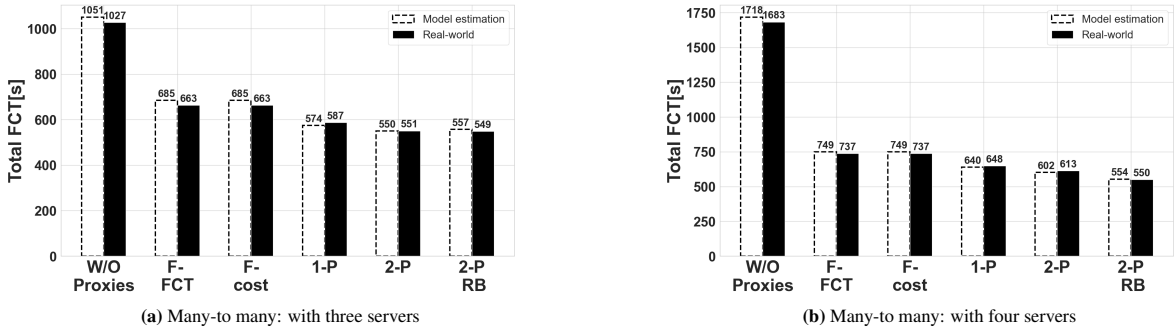


Figure 15: Real-world experiment: We measure the total FCT obtained in a real-world experiment on GCP. We consider a many-to-many topology with users in Brazil, England, Finland, and Japan. We use our CloudPilot system to deploy cloud proxies with Kubernetes and HAProxy, such that our algorithms select the locations. We then compare the results against our model prediction. We can see that our FCT acceleration prediction achieves very close results to the obtained real-world results.

with 60 flows in unoptimized Python, as in Fig. 7(d), takes less than one second for the flow-greedy algorithms while for the 1-P, 2P, and 2P with RB, it takes 1.6s, 6.5s, and 131s, respectively.

7.3. Cloud Experiments

Methodology. We deploy our CloudPilot system on GCP to set up the proxies and run iPerf3 tests (details are provided in §4). We consider a many-to-many full-mesh use-case for 2 cases: (1) with three servers in Hamina (Finland), Sao-Paulo (Brazil), and Tokyo (Japan), and therefore six flows. (2) with four servers in Hamina (Finland), London (England), Sao-Paulo (Brazil), and Tokyo (Japan), and therefore twelve flows. In both cases, we measure the FCT of each flow. All the hosts run virtual machines with default instances (E2-medium). The budget is 50¢ per flow, *i.e.*, is 3\$ for the first experiment and 6\$ for the second one. Each result averages 20 runs.

Results. Fig. 15 shows how the proxy-greedy algorithms achieve better real-world results than the flow-greedy ones, as previously seen in the simulations. This is because the proxy-greedy algorithms family tends to

prefer the one-proxy connection method. Significantly, as we also saw in the model evaluation for individual flows (Fig. 4), our modeled predictions for the total system FCT appear close to the real-world measured FCTs. In addition, we can see that the FCT acceleration can increase when we have a larger budget and more flows with long RTT.

8. Conclusion and future work

In this paper, we introduced *CloudPilot*, a Kubernetes-based system that measures communication parameters across different cloud regions, and uses these measurements to deploy cloud proxies in optimized locations on multiple cloud providers. We further demonstrated how it can significantly improve the flow completion time of global geo-distributed applications by relying on an optimized placement of cloud proxies.

Many ideas and questions remain open and left for future work. There is a vast area for improvement that can be made in the CloudPilot system. For instance, we want to perform cloud evaluation in other cloud providers such as AWS or Azure. In addition, with the

rise of multi-cloud applications, it will be interesting to check the multi-cloud scenario, when proxies can be allocated in more than one cloud provider, with the potential to reach better results.

Also, we could support other cloud providers like Azure, or use Kubernetes container platforms that support multi-cloud like OpenShift [37] or Google Anthos [12].

Going forward it will be interesting to examine our problem in the context of Service Function Chaining when we do not use just TCP split proxy acceleration, but integrate further acceleration and services, such as compression, caching, transcoding to different QoS levels and their combination with encryption, firewalls, and other network services.

Another area of exploration concerns sharing proxies among multiple flows and considering additional cost parameters, such as the cost of egress traffic in a cross-cloud setup in addition to the cost of proxies. In addition, the use of the presented approach for edge cloud communication is an important potential direction for further exploration. For example, uploading large video files for backup from security cameras to the cloud can benefit significantly from our CloudPilot system and save uploading time. Likewise, we plan to explore whether CloudPilot can be extended to handle large transfers with different priorities and QoS in complex cloud edge applications, such as AR/VR, manufacturing, digital twin, etc.

Acknowledgment

The authors would like to thank Ofer Biran, Roy Mitrany and Aran Bergman for useful discussions. This work was partly supported by the Louis and Miriam Benjamin Chair in Computer-Communication Networks, the Israel Science Foundation (grant No. 1119/19), and the Hasso Plattner Institute Research School.

References

- [1] Bergman, A., Cidon, I., Keslassy, I., Rotman, N., Schapira, M., Markuze, A., Zohar, E., 2018. Pied piper: Rethinking internet data delivery. arXiv preprint arXiv:1812.05582 .
- [2] Bhattacharjee, D., Tirmazi, M., Singla, A., 2017. A cloud-based content gathering network, in: USENIX HotCloud, Santa Clara, CA.
- [3] Cai, C.X., Le, F., Sun, X., Xie, G.G., Jamjoom, H., Campbell, R.H., 2016. CRONets: Cloud-routed overlay networks, in: IEEE ICDCS. doi:10.1109/ICDCS.2016.49.
- [4] Cardwell, N., Cheng, Y., Gunn, C.S., Yeganeh, S.H., Jacobson, V., 2016. Bbr: Congestion-based congestion control: Measuring bottleneck bandwidth and round-trip propagation time. Queue 14, 20–53.
- [5] Cohen, E., Krishnamurthy, B., Rexford, J., 1998. Improving end-to-end performance of the web using server volumes and proxy filters, in: ACM SIGCOMM, New York, NY, USA.
- [6] Dang, T.K., Mohan, Nitinder and Corneo, Lorenzo and Zavadovski, Aleksandr and Ott, Jörg and Kangasharju, Jussi, 2021. Cloudy with a chance of short RTTs: analyzing cloud connectivity in the internet, in: ACM IMC, pp. 62–79.
- [7] Deng, J., Tyson, G., Cuadrado, F., Uhlig, S., 2017. Internet scale user-generated live video streaming: The Twitch case, in: Passive and Active Measurement, pp. 60–71.
- [8] Farkas, V., Héder, B., Nováczki, S., 2012. A split connection TCP Proxy in LTE Networks, in: Information and Comm. Technologies.
- [9] Fischlin, M., Günther, F., 2014. Multi-stage key exchange and the case of Google’s QUIC protocol, in: Proceedings of the 2014 ACM SIGSAC Conference on Computer and Communications Security, pp. 1193–1204.
- [10] Ge, C., Wang, N., Chai, W.K., Hellwagner, H., 2018. QoE-assured 4K HTTP live streaming via transient segment holding at mobile edge. IEEE J. Select. Areas Commun. .
- [11] Google, 2022. GCP inter region latency. URL: <https://datastudio.google.com/reporting/1c733b10-9744-4a72-a502-92290f608571/page/70YCB>.
- [12] Google Cloud, 2022a. Anthos technical overview. URL: <https://cloud.google.com/anthos/docs/concepts/overview>.
- [13] Google Cloud, 2022b. General-purpose machine family. URL: <https://cloud.google.com/compute/docs/general-purpose-machines>.
- [14] Google Cloud, 2022c. Price list (GCP VM pricing). URL: <https://cloud.google.com/pricing/list>.
- [15] Google Cloud, 2022d. Virtual private cloud pricing (GCP network pricing). URL: <https://cloud.google.com/vpc/pricing>.
- [16] Guo, Y., Ge, Z., Urgaonkar, B., Shenoy, P., Towsley, D., 2004. Dynamic cache reconfiguration strategies for a cluster-based streaming proxy, in: Web Content Caching and Distribution.
- [17] Haq, O., Doucette, C., Byers, J.W., Dogar, F.R., 2020. Judicious QoS using cloud overlays, in: ACM CoNEXT, pp. 371–385.
- [18] iPerf3, 2022. iPerf - the ultimate speed test tool. URL: <https://iperf.fr/iperf-download.php>. accessed: 2022-30-01.
- [19] Jonathan, A., Ryden, M., Oh, K., Chandra, A., Weissman, J., 2017. Nebula: Distributed edge cloud for data intensive computing. IEEE TPDS .
- [20] Kelly, F., 2001. Mathematical modelling of the internet, in: Mathematics unlimited—2001 and beyond. Springer, pp. 685–702.
- [21] Kim, B.H., Calin, D., Lee, I., 2017. Enhanced split TCP with end-to-end protocol semantics over wireless networks, in: IEEE WCNC.
- [22] Kloudas, K., et al., 2015. Pixida: Optimizing data parallel jobs in wide-area data analytics. Proc. VLDB Endow. .
- [23] Kopparty, S., Krishnamurthy, S., Faloutsos, M., Tripathi, S., 2002. Split TCP for mobile ad hoc networks, in: IEEE Globecom. doi:10.1109/GLOCOM.2002.1188057.
- [24] Lai, F., Chowdhury, M., Madhyastha, H., 2018. To relay or not to relay for Inter-Cloud transfers?, in: USENIX HotCloud, Boston, MA.
- [25] Lam, C.F., 2021. (invited) Google Fiber Deployments: Lessons learned and future directions, in: OFC.
- [26] Lam, C.F., Liu, H., Koley, B., Zhao, X., Kamalov, V., Gill, V., 2010. Fiber optic communication technologies: What’s needed for datacenter network operations. IEEE Communications Mag-

azine Vol.48 No.7.

- [27] Le, F., Nahum, E., Kandlur, D., 2016. Understanding the performance and bottlenecks of cloud-routed overlay networks: A case study, in: ACM Workshop on Cloud-Assisted Networking, p. 7–12.
- [28] Lepikhov, K., 2022. Propagation delay in Géant. URL: <https://wiki.geant.org/display/public/EK/PropagationDelay>.
- [29] Li, P., Guo, S., Miyazaki, T., Liao, X., Jin, H., Zomaya, A.Y., Wang, K., 2017. Traffic-aware geo-distributed big data analytics with predictable job completion time. IEEE TPDS .
- [30] Luglio, M., Sanadidi, M., Gerla, M., Stepanek, J., 2004. On-board satellite "split TCP" proxy. IEEE J. Select. Areas Commun. .
- [31] Markuze, A., Bergman, A., Dar, C., Keslassy, I., Cidon, I., 2020. Kernels of splitting TCP in the clouds, in: Netdev 0x14.
- [32] Martello, S., Toth, P., 1980. Solution of the zero-one multiple knapsack problem. European J. of Op. Research .
- [33] Mok, R.K., Zou, H., Yang, R., Koch, T., Katz-Bassett, E., Claffy, K., 2021. Measuring the network performance of Google Cloud Platform, in: ACM IMC, pp. 54–61.
- [34] Pathak, A., ang, Y.A., Huang, C., Greenberg, A., Hu, Y.C., Kern, R., Li, J., Ross, K.W., 2010. Measuring and evaluating TCP splitting for cloud services, in: PAM'10, Zurich, Switzerland.
- [35] Pierre, G., van Steen, M., 2006. Globule: a collaborative content delivery network. IEEE Communications Magazine .
- [36] Pu, Q., Ananthanarayanan, G., Bodik, P., Kandula, S., Akella, A., Bahl, P., Stoica, I., 2015. Low latency geo-distributed data analytics. ACM SIGCOMM CCR .
- [37] RedHat, 2022. OpenShift platform: white paper. URL: <https://www.mindtree.com/insights/resources/redhat-openshift-container-platform>.
- [38] Roşu, M.C., Roşu, D., 2002. An evaluation of TCP splice benefits in web proxy servers, in: ACM WWW, p. 13–24.
- [39] Rotman, N.H., Ben-Itzhak, Y., Bergman, A., Cidon, I., Golikov, I., Markuze, A., Zohar, E., 2022. Cloudcast: Characterizing public clouds connectivity. arXiv preprint arXiv:2201.06989 .
- [40] Shalev, L., Ayoub, H., Bshara, N., Sabbag, E., 2020. A cloud-optimized transport protocol for elastic and scalable HPC. IEEE Micro 40, 67–73. doi:10.1109/MM.2020.3016891.
- [41] SpeedTest, 2022. Public SpeedTest servers. URL: <https://as62240.net/speedtest>.
- [42] Taft, R., , Sharif, I., Matei, A., VanBenschoten, N., Lewis, J., Grieger, T., Niemi, K., Woods, A., Birzin, A., Poss, R., 2020. Cockroachdb: The resilient geo-distributed SQL database, in: ACM SIGMOD.
- [43] Toledo, K., et al., 2022. CloudPilot system git. URL: https://github.com/kfirtoledo/CloudPilot_Project.
- [44] Wu, J., Ravindran, K., 2009. Optimization algorithms for proxy server placement in content distribution networks, in: IFIP/IEEE International Symposium on Integrated Network Management-Workshops.
- [45] Yeganeh, B., Durairajan, R., Rejaie, R., Willinger, W., 2020. A first comparative characterization of multi-cloud connectivity in today's internet, in: Passive and Active Meas., pp. 193–210.
- [46] Zhang, F., Fu, X., Yahyapour, R., 2017. CBase: A new paradigm for fast virtual machine migration across data centers, in: IEEE/ACM CCGRID.
- [47] Zhang, H., et al., 2019. Harmony: An approach for geo-distributed processing of big-data applications, in: IEEE CLUSTER.