

HybridNN: Supporting Network Location Service on Generalized Delay Metrics

Yongquan Fu, Yijie Wang, and Ernst Biersack,

Abstract—Distributed Nearest Neighbor Search (DNNS) locates service nodes that have shortest interactive delay towards requesting hosts. DNNS provides an important service for large-scale latency sensitive networked applications, such as VoIP, online network games, or interactive network services on the cloud. Existing work assumes the delay to be symmetric, which does not generalize to applications that are sensitive to one-way delays, such as the multimedia video delivery from the servers to the hosts. We propose a relaxed inframetric model for the network delay space that does not assume the triangle inequality and delay symmetry to hold. We prove that the DNNS requests can be completed efficiently if the delay space exhibits modest inframetric dimensions, which we can observe empirically. Finally, we propose a DNNS method named HybridNN (*Hybrid Nearest Neighbor search*) based on the inframetric model for fast and accurate DNNS. For DNNS requests, HybridNN chooses closest neighbors accurately via the inframetric modelling, and scalably by combining delay predictions with direct probes to a pruned set of neighbors. Simulation results show that HybridNN locates nearly optimally the nearest neighbor. Experiments on PlanetLab show that HybridNN can provide accurate nearest neighbors that are close to optimal with modest query overhead and maintenance traffic.

I. INTRODUCTION

Latency-sensitive applications, such as P2P based VoIP and IPTV [1], interactive network services on the cloud (e.g., Office Live Workspace [2], Google Maps [3]), online network games, need to transmit data from geo-distributed servers (called a service node) in real-time to many hosts quickly. High transmission delays reduce the Quality of Experience (QoE) of users [4], which lead to significant business losses [5]. For instance, Google reports that its revenue decreases by 20% when the latency of showing search results increases by 500 ms; similarly, Amazon claims that its sales amount decreases by 1% if the page-response latency increases by 100 ms [5].

Since there are hundreds or thousands of service nodes that provide identical services to hosts, there is an increasing push for service providers to route real-time data to a host from geo-distributed servers that are nearest to that host. For example, Google routes users' search queries to geographical-nearby servers [6]; Akamai redirects hosts' content requests to replica servers mainly based on proximity conditions [7]; CoralCDN [8] uses OASIS [9] and DONAR [10] to select proxy servers near to end hosts based on geographic distances. However, selecting nearest servers to hosts are still far from standard due to several challenges.

Yongquan Fu and Yijie Wang are with National Key Laboratory for Parallel and Distributed Processing, College of Computer Science, University of Defense Technology.

Ernst Biersack is with Networking and Security Department, Eurecom.

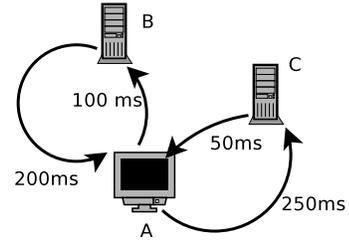


Fig. 1. Illustrating the RTT and OWDs. Suppose B and C are two servers that are able to supply short videos to host A . If we use the RTT metric to minimize the delay of video delivery, we may arbitrarily choose any of them to send videos to host A based on the RTT metric, since the RTT between A , B and that between A , C are all 300 ms. However, since the video files are transmitted from servers to hosts, the OWDs from servers to hosts become more important [16]. We can see that the OWD from server C to host A is four times less than that from server B to host A . Therefore, choosing server C to serve host A significantly minimizes the content transmission delay for host A , which is feasible only when we use the OWD metric for delay optimizations.

First, *selecting nearest servers must prove to be reliable, since service providers need to ensure the QoE fairly for all hosts*. Selecting nearest servers using proximity coordinates [11], [12] or geographic distances [9] suffer from the mismatch between the estimated delays and real-world delays [6], which makes the selection accuracy hard to be predicted. On the other hand, selecting nearest servers using distributed search such as Meridian [13] or OASIS [9] avoid such mismatch problems using direct probes, but may terminate at service nodes that are much worse than the nearest ones, since the search is easily trapped into local minima due to the clustering [14] and Triangle Inequality Violations (TIV) [15] properties of the delay space.

Second, *selecting nearest servers must be aware of unidirectional delays whenever possible*. Since routing on the Internet is asymmetric [16], the delays from servers to hosts may deviate those in the reverse direction in several times. Furthermore, One-Way Delay (OWD) measurements become increasingly practical due to the advance of measurement techniques such as OWAMP [17] or Reverse Traceroute [18]. However, delay optimizations using Round Trip Time (RTT) ignores such delay asymmetry. For multimedia streaming, application-level multicast, or more generalized applications where data flows in one directions, such agnostics of unidirectional delays degrades the effectiveness of selected servers, as shown in Fig 1.

Third, *selecting nearest servers must find good tradeoff between the response time and timeliness*. The response time lasts several seconds for server selections using on-demand

arXiv:1108.1928v1 [cs.DC] 9 Aug 2011

probing such as Meridian [13] or OASIS [9]. However, the response time degrades the QoE of users in latency-sensitive applications, such as online workspace, online music. OASIS caches nearest servers for each IP prefix using in-advance probes once a week, which has better response time. However, the cached server selections tend to be suboptimal, since the delays vary due to routing dynamics or server workloads [19], and service nodes may be added or removed dynamically. Therefore, it is difficult to find good tradeoff between response time and the timeliness of server selections.

The goal of this paper is to provide new algorithms to address the first two challenges. To this end, we develop a general enough delay model that captures the major statistics of the delay space, including: TIV, delay dynamics and asymmetry of delays. This paper makes three contributions.

First, we analytically demonstrate that we can find approximately nearest servers quickly by iteratively searching closer nodes to the host using sampled nodes from proximity regions of each node. However, the analytical method requires a large number of samples, which does not scale well.

Second, we introduce a novel distributed algorithm, named HybridNN, that finds nearest service nodes for any machine on the Internet (called a target). This algorithm derives from our analytical method, which preserves the accuracy and speediness of the analytical method. However, HybridNN has better dynamic adaptation and reduced measurement costs.

(i) **Dynamic adaptation.** A practical DNNS algorithm needs to proactively maintain moderate service nodes as samples for DNNS queries, irrespective of the system dynamics. HybridNN dynamically maintains such neighbors using a concentric ring used in Meridian [13] or OASIS [9]. However, HybridNN has two improvements:

- The maximum number of nodes stored per ring is derived from the lower bounds of required samples in the analytical method, which implies that HybridNN requires the lowest possible number of samples that has the same accuracy guarantee as the analytical method.
- HybridNN proposes a biased sampling based concentric ring maintenance scheme, in order to sample enough nodes for each ring. Specifically, different from previous neighbor discoveries based on a gossip protocol, we also periodically discover a small number of nearest nodes and farthest nodes to each node as neighbors in the concentric ring. This is because given a concentric ring, the innermost and outermost rings contain only a few neighbors compared to other rings, which are hardly to be sampled using a gossip based neighbor discovery protocol.

(ii) **Reducing measurement costs.** HybridNN adopts scalable delay predictions to reduce the measurement costs.

- HybridNN maintains the concentric rings using estimated pairwise delays with the revision [20] of the Vivaldi network coordinate [21], which significantly reduces the maintenance overhead of HybridNN compared to Meridian.
- HybridNN selects candidate neighbors that are close to the target using delay predictions. Since delay predictions

are only approximations of real-world delays, HybridNN also uses a small number of delay probes to avoid being misled by inaccurate delay predictions. Interestingly, although the network coordinate distances are symmetric, we empirically find that our hybrid delay measurement approach provides the accurate nearest next-hop neighbor for both symmetric and asymmetric delay data sets. This is because we replace inaccurate coordinate distances with direct probes using the error indicator of Vivaldi coordinate, which relieves the mismatch between symmetric coordinate distances and asymmetric delays.

Third, we validate our algorithm using real-world delay data sets and PlanetLab deployments. Through simulation study, we show that HybridNN finds servers close to optimal for symmetric and asymmetric delay data sets. In fact, in more than 95% of cases, HybridNN locates the ground-truth nearest servers for the targets. Furthermore, most queries terminate within four search hops, which implies that HybridNN can return the search results fast. Using PlanetLab deployments, we confirm that HybridNN can locate accurate nearest servers with low query loads and control overhead, with moderate query time that improves Meridian in more than 15% of cases.

II. SYSTEM MODEL

A. Problem Definition

In this section, we formally define the nearest server location problem. Let V denote a set of service nodes and hosts. Let a distance function d denote the pairwise delays between node pairs in V . Let N be the number of service nodes.

Our objective is to minimize the serving delays of latency-sensitive applications by finding a service node for a requesting host with the minimum delay. As discussed in the previous section, we expect a generalized delay optimization scenario where the delay may be symmetric or asymmetric according to the problem context and measurement tools. Furthermore, the service nodes may be added or removed, which causes system churns. As a result, we need to locate the service node that is closest to the target from dynamic service nodes.

We study a distributed approach to realize our objective, since the centralized approach has several well-known weaknesses, including: it requires global delay measurements that is hard to obtain for dynamic service nodes; it incurs the single point of failures. On the other hand, the distributed approach avoids such weaknesses through collaborations of service nodes. Specifically, we formulate the **Distributed Nearest Neighbor Search** (DNNS) as:

Definition II.1. (*Distributed Nearest Neighbor Search*): For a set of dynamic service nodes, given any target T on the Internet, the objective of the Distributed Nearest Neighbor Search is to find one service node that has the smallest delay to T , based on the distributed collaboration of service nodes.

The definition of DNNS is not novel, since existing research on closest server discovery [22], [23], [12], [13], [9], [10] has formulated the similar problem. Intuitively, DNNS consists of multiple steps. At each step, a current service node P tries to

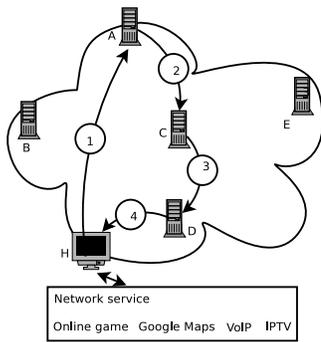


Fig. 2. A DNNs query service substrate for network services.

locate a new service node that is closer to the target T than node P . The flowchart of a sample DNNs query is shown in Fig 2. When a host T accesses a networked service, the local service client module creates a DNNs query to locate the nearest service machine to the client T . The query message is firstly forwarded to the bootstrap machine of the DNNs service (Step 1). Then our DNNs query system will forward the query message recursively until locating a nearest service machine (Step 2 \rightarrow 3). Finally, our system returns the contact addresses of the found service nodes to host T (Step 4).

B. Key DNNs Requirements

To be useful for latency-sensitive applications, we identify key goals for the DNNs:

- **Accurate**, we need to find a service node with the lowest interactive time in order to increase the Quality of Experience of users.
- **Fast**, we need to obtain the nearest service node with low query periods. Otherwise, long query time makes the DNNs less attractive for server redirections in latency-sensitive applications.
- **Scalable**, the DNNs process should incur low bandwidth costs with increasing system size.
- **Resilient to churns**, the DNNs process should find accurate results when the service nodes crash or new service nodes are added.

C. Discussion

Since the DNNs process may last several seconds due to on-demand probing, performing DNNs for each query from hosts may even hurt the Quality of Experience of users, which is significant for small Web objects. For example, Google typically returns responses in less than 0.4 seconds; however, such low response periods are difficult to be realized when applying the DNNs process before returning the responses.

Therefore, in order to realize a practical nearest server redirection service, we need to proactively run DNNs for each host and redirect hosts' requests using cached DNNs results, in order to achieve millisecond-level response time. For example, OASIS [9] shows that it is feasible to cache DNNs queries of IP prefixes for server redirections without reducing the DNNs accuracy.

We do not study how to organize cache results in this paper; instead, we assume that a DNNs caching service exists to map hosts' requests to nearest servers using cached DNNs queries. Our focus is to realize an accurate, scalable and resilient DNNs system with low DNNs query periods. Since if the DNNs query last long periods, then crawling DNNs for every IP prefix will be less efficient.

III. RELATED WORK

First, for the theoretical computer science field, research on the nearest neighbor search mainly focuses on designing efficient algorithms in the metric space [24], [25], [26], [27]. However, applying algorithms in the metric space into DNNs is inappropriate, since the delay space violates the triangle inequality that is required by the metric space model [20].

On the other hand, for the network system field, research on nearest neighbor search can be classified into centralized and distributed approaches according to the communication patterns of the search process.

A. Centralized Approaches

The centralized scheme uses a centralized sorting process to select nearest neighbors for target nodes. However, the centralized approach does not scale well with increasing system size, since collecting and transmitting the distance measurements easily cause performance bottlenecks, which degrades the service availability.

Guyton et al. [11] pioneer the research on finding the closest server replica in a centralized manner. They use the Hotz's metric [28] to represent pairwise hop distances using $O(N)$ measurements to landmark nodes, where N denotes the number of server replicas. However, smaller hop distances do not mean the shorter delays, because one hop may pass continents or a data center. Later Carter and Crovella [29], [30] combine the RTT and available bandwidth measurements to dynamically select optimal server replica with minimal response time. However, the dynamic server selection approach does not scale well due to the quadric measurement costs. Netvigitor [31] collects RTT values from hosts to landmarks and milestone nodes based on the Traceroute measurements, and estimates nearest servers based on local clustering. However, Netvigitor does not guarantee the estimation accuracy, and may get obsolete results since Netvigitor does not perform active measurements. Different from Netvigitor, CRP [32] leverage the dynamic association of nodes with replica servers from CDNs to determine the proximity between end hosts. CRP incurs low maintenance costs similar as Netvigitor. However, CRP does not guarantee the accuracy. iPlane [33], [34] constructs a synthetic topology structure for the Internet. iPlane estimates the nearest servers using the approximated delays on the synthetic topology. However, in order to provide services for hosts spanning geo-distributed places, iPlane consumes heavy bandwidth costs to perform active measurements.

B. Distributed Approaches

The DNNs approach iteratively selects closer nodes using distributed nearest neighbor search by local measurements

towards a small set of neighbors, which reduces the network measurement overhead and is more scalable than the centralized approach. Existing DNNS methods fall into four families based on their search rules: (i) Bin based DNNS; (ii) Topology based DNNS; (iii) Greedy search based DNNS; (iv) Ring search based DNNS.

Bin based DNNS. Ratnasamy et al. [22] assign nodes into "bins" based on the ordered sequence of RTT measurements to landmarks, and declare nodes are close to each other in the same bin. However, the bin approach does not guarantee the accuracy, and fails when the landmarks crashes.

Topology based DNNS. Tiers [35] locates the nearest nodes by a top-down approach with a hierarchical clustering tree, but may cause load imbalance for nodes near the root of the tree. Besides, Tiers do not guarantee the search accuracy since the tree does not strictly preserve the pairwise proximity.

Greedy search based DNNS. Mithos [23] iteratively locates proximate neighbors with $O(N)$ hops by a gradient descent based protocol in the overlay construction, but terminates earlier before locating the real nearest nodes due to the limited diversity in the neighbor set. PIC [12] iteratively locates nearest neighbors at each search step in terms of the coordinate distance. However, PIC is prone to be trapped into the local minima since the coordinate distance only approximates the delays. DONAR [10] redirects host requests to optimal server replicas by considering the network proximity, the routing optimization and server loads. DONAR uses geographic distances as the proximity metric in order to reduce measurement costs. However, DONAR may find suboptimal server replicas for delay minimizations since the delay values are not consistent with the geographic distances.

Ring search based DNNS. Our work is closely related to Meridian [13], which seeks approximately nearest nodes in $\log(N)$ steps. Meridian [13] maintains a loosely connected overlay using a gossip based peer finding scheme. The neighbors are organized in concentric rings with exponentially increasing radii. For a DNNS request, Meridian iteratively locates one next-hop node that is β ($\beta < 1$) times closer to the target T than the current Meridian node. Compared to other families of DNNS, Meridian is more accurate by using rings of neighbors that promote the diversity of neighbor sets [13]. However, several studies have identified that Meridian may fail to find the closest service node due to the last-hop clustering of servers [14], and TIV of the network delay space [20]. Similar as Meridian, OASIS [9] organize neighbors as concentric rings for each service node, and iteratively search nearest service node for the request host in terms of the geographic distances. OASIS reduces the delay measurement costs in Meridian through the static geographic coordinates, and has low response time using in-advance probes. However, OASIS does not guarantee the accuracy of the search results, since selecting the geographically closest servers may incur high delays [6].

To address these problems, two adjustments are proposed: (i) explicitly finding the clustering subsets based on the structure of IP addresses [14] or, (ii) adding additional neighbors for DNNS that may not be chosen due to the TIVs [20]. However, finding the clusters of nodes sharing identical last

hops becomes insufficient when the service nodes spread over nearby subnets, which may still mislead the DNNS queries due to no forwarding nodes closer enough to the target. Furthermore, finding all neighbors that are affected by the TIVs is challenging since calculating the TIVs for decentralized service nodes is very difficult; besides, adding additional neighbors for DNNS also increases the query overhead. Due to the limitations of modifications for Meridian, significant challenges remain in DNNS. We focus on tackling these challenges in this paper.

IV. DATA SETS

Our empirical data sets include four publicly available real-world RTT data sets, covering the delay measurements between wide-area DNS servers and those between end hosts [36]. (i) **DNS3997.** A RTT matrix collected between 3997 DNS servers by Zhang et al. [37] using the King method [38]. The matrix is symmetric in that $d_{ij} = d_{ji}$, for any pair of items i and j , where d denotes the delay matrix. (ii) **Host479.** A RTT delay matrix based on RTT measurements that last 15-day periods between the Vuze BitTorrent clients [39]. Host479 is asymmetric, where in over 40% of the cases delay pairs d_{AB} and d_{BA} in Host479 differ more than 4 times. This is because RTT measurements between node pairs are not synchronized and delay results are affected by varying queueing delays at end hosts [39]. (iii) **DNS1143.** A RTT matrix between 1143 DNS servers collected by the MIT P2PSim project [40] using the King method [38]. The matrix is symmetric in that $d_{ij} = d_{ji}$, for any pair of items i and j , where d denotes the delay matrix. (iv) **DNS2500.** A RTT matrix between 2500 DNS servers by the Meridian project [13] using the King method. The matrix is also symmetric.

Since obtaining the one-way delays between large-scale nodes is extremely difficult, we use Host479 as an asymmetric delay data set. However, we do not claim that our experiments on Host479 are the same as those on the one-way delay metric.

V. A GENERALIZED DELAY MODEL FOR THE DELAY SPACE

In this section, we present a simple and general enough delay model for the delay space. Our model captures the important characteristics of the delays, including TIV, dynamics and asymmetry of RTTs and OWDs. In the next section, we will analyze the DNNS problem on our model.

Assuming that we select a node P in V as the center of a ball, and choose a positive real number r as the radius of the ball, then we call a **closed ball** $B_P(r)$ as the set of nodes whose delays to node P are not larger than r , i.e., $B_P(r) = \{v | d(P, v) \leq r, P, v \in V\}$. Furthermore, the **volume** of a ball is the number of nodes covered by the ball. Besides, we define the cover relation of different set of nodes as follows:

Definition V.1 (Cover). *Let S and Ω be two sets of nodes, if $\Omega \subseteq S$, then the set S is said to cover the set Ω .*

A. Definition

We first state the requirements for a delay model suitable for RTTs and OWDs used for delay minimizations. (i) The

delay model should relax the symmetry requirements, since the OWDs are asymmetric due to routing asymmetry [41]. Besides, although RTT is symmetric by only accounting for the delays on the routing paths, real-world RTT measurements may be asymmetric due to variations of queuing delays at end hosts or un-synchronized measurements [42]. (ii) The delay model d should allow TIV to exist, since the RTT metric exhibits TIV [15]. (iii) The delay model d should allow dynamic delays, since the delay varies from time to time [19].

Therefore, inspired by the inframetric model [43] that allows the TIVs, we extend the inframetric model to a relaxed inframetric model that relaxes the symmetry requirement, where the distance function d satisfies:

Definition V.2 (Relaxed Inframetric Model). *Let a distance function $d : V \times V \rightarrow \mathbb{R}^+$ be a relaxed ρ -inframetric ($\rho > 1$), if d satisfies the following conditions for any pair of nodes u and v : (1) if $d(u, v) = 0$, then $u = v$; (2) $d(u, v) \leq \rho \max\{d(u, w), d(v, w)\}$, for any arbitrary node w satisfying $w \notin \{u, v\}$.*

Pros of the Relaxed Inframetric Model: The condition (2) in Def V.2 states a generalized relation of any directed triple from V , which has two beneficial properties:

- **TIV-adaptive.** Intuitively, smaller ρ implies that three edges are closer to each other; while larger ρ implies that one edge is significantly larger than any of the other two edges, which may introduce a TIV. Therefore, similar as the inframetric model, the relaxed inframetric model naturally allows the occurrence of TIVs.
- **Dynamics-adaptive.** The inframetric model allows the delay variations by varying the inframetric parameter ρ to describe the relations of updated triples. Therefore, both inframetric model and the relaxed inframetric model are able to model variations of triples due to delay variations.
- **Asymmetry-aware.** The relaxed inframetric model allows the asymmetry in the delay space, which generalizes to RTTs and OWDs. As a result, we are able to analyze DNNS on symmetric and asymmetric delays through the relaxed inframetric model.

Having shown the advantages of the relaxed inframetric model, next we discuss the statistical property of the inframetric parameter ρ .

First, the seminal work states that if the delay space obeys the triangle inequality, then ρ must be smaller or equal than 2 [43]. However, when ρ is smaller than 2, there may exist TIVs. For example, given a triple with pairwise RTTs 3, 1, 1.8, we can see that the inframeter parameter ρ is approximately 1.67 but there also exists a TIV in the triple. Therefore, we can see that $\rho \leq 2$ is only a *necessary* but not a *sufficient* condition for no TIVs.

Second, we find that the inframetric parameter ρ is quite low for most triples. First, the 95th percentiles of all data sets of ρ are below 2.5. Low inframetric parameter ρ means the largest edges in triples are not too much larger than the other edges of the triples. Second, among the triples whose ρ are bigger than 2, their ρ values are around 3 on average. Therefore, selecting $\rho=3$ is reasonable to model most of the triples.

B. Dimensions on the Relaxed Inframetric Model

Having introduced the definition of the relaxed inframetric model, now we analyze the growth dimension of the relaxed inframetric model, which is the ratio of the number of nodes covered by two closed balls with the identical center and varying radii [43], [44].

The growth dimension is important for efficient DNNS. As shown by Karger and Ruhl [44], assuming that the growth dimension is low, each node P can uniformly sample a modest number of nodes to locate a node that is closer to any other node in V . Therefore, we can recursively find nodes closer to the target based on the above sampling procedure, which helps the design of the DNNS algorithms. However, since Karger and Ruhl assumes the triangle inequality to hold [44], we cannot immediately apply their DNNS results into the relaxed inframetric model. Accordingly, we need new proof techniques for DNNS analysis.

The growth dimension for the inframetric space [43] is defined as follows:

Definition V.3 (Growth [43]). *For a ρ -inframetric model, for any $r \in \mathbb{R}^+$ and $P \in V$, if $|B_P(\rho r)| \leq \gamma_g |B_P(r)|$, where $\gamma_g \in \mathbb{R}^+$, the ρ -inframetric model is said to have a growth $\gamma_g \geq 1$.*

The growth dimension γ_g on the inframetric model generalizes the growth definition in the metric space which assumes the triangle inequality to hold [44], [37]. Therefore, the growth γ_g inherits the intuitive meanings of the growth definition in the metric space. Specifically, low growth γ_g means that the number of nodes covered by the closed ball $B_P(\rho r)$ is comparable to the number of nodes covered by the closed ball $B_P(r)$. Therefore, when we expand a ball around a node $P \in V$, we can see that new nodes in V "come into view" at a constant rate [44].

Finally, based on Def V.3, the infimum of the growth dimension γ_g equals the ratio of the volume between $B_P(\rho r)$ and $B_P(r)$ for any node P and radius r . Since we are interested in the infimum, when we refer to the growth of the inframetric space, we mean the infimum accordingly.

Next, we empirically evaluate the growth dimension of the delay space with respect to the radius r and the inframetric parameter ρ . Our evaluation complements the seminal work on the growth in the inframetric model [43] using symmetric and asymmetric data sets. Recall that computing the growth is trivial by comparing the volumes of the balls with identical centers and varying radii.

Fig 3 shows the median and 90th percentile growth values for varying radii. The median growth of most data sets is relatively small, and declines quickly with increasing radii for most data sets except for Host479. For Host479, the median growth may increase as the radii increase. On the other hand, the 90th percentile growth shows divergent dynamics for different data sets, revealing "M"-shape dynamics, indicating that a small fraction of growth values may increase or decrease with increasing radii.

Furthermore, by selecting different percentages of nodes for the statistics, Fig 3 shows that the median growth is less sensitive to the sample size compared to the magnitudes of

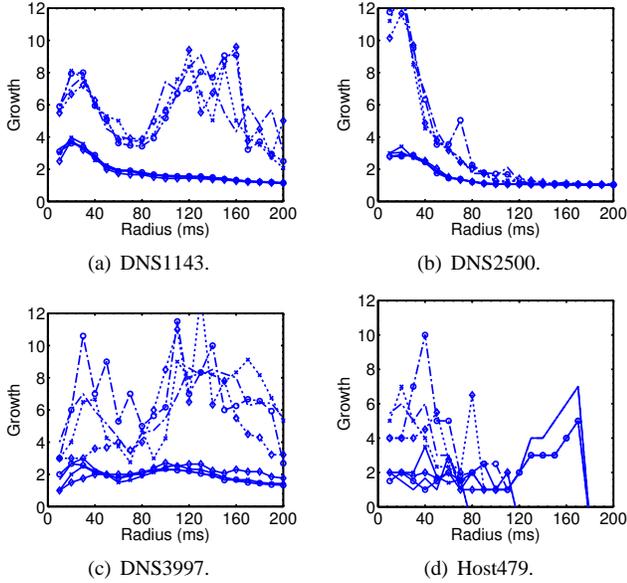


Fig. 3. The statistics of the median and 90-th percentile growth γ_g for $\rho = 3$; $-\diamond-$ denotes median values computed from sampled 20% nodes; $-x-$ denotes median values computed from sampled 50% nodes; $-o-$ denotes median values computed from sampled 75% nodes; $-$ represents median values computed from all nodes; $\dots\diamond\dots$ denotes 90-percentile values computed from sampled 20% nodes; $\dots x \dots$ denotes 90-percentile values computed from sampled 50% nodes; $\dots o \dots$ denotes 90-percentile values computed from sampled 75% nodes; $-.-$ represents 90-percentile values computed from all nodes.

radii; while the 90th percentile growth becomes relatively more sensitive to the sample size.

In summary, the growth metric γ_g of the delay space is quite low. Furthermore, with increasing radius, the growth γ_g decreases to 2 quickly on average. However, sometimes the growth values increase for increasing radius, which means that there are many nodes that have similar distances to each other. This usually corresponds to cases where the center of the ball is a node on the edge of a cluster, where nodes in the same cluster have smaller distances compared to those to other nodes not in the same cluster.

VI. EFFICIENT DNNS ON THE RELAXED INFRAMETRIC MODEL

In this section, using the relaxed inframetric model presented in Sec V, we analyze how to design an efficient DNNS using localized operations suitable for distributed systems. Proofs are omitted due to space limits, which can be found in the full report [36].

Our major result is that it is feasible to design an accurate and fast DNNS algorithm for the relaxed inframetric mode, at the expense of sampling enough candidate servers from the proximity region of each node. We construct a simple DNNS process satisfying our major result. However, the simple DNNS process incurs relatively high measurement costs due to the sampling conditions, which will be improved in the next section.

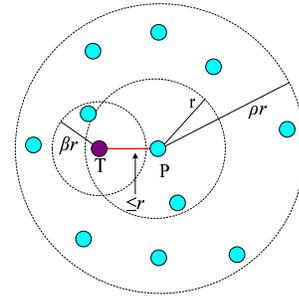


Fig. 4. Sampling closer nodes to a target T from $B_P(\rho r)$ in the ρ -inframetric model with growth γ_g .

A. Sampling Conditions to Locate Closer Nodes To Targets

In this section, We analyze samples required to locate a node closer to a target based on the growth dimension in Sec V-B. The sampling conditions serves as the basis for the efficient DNNS algorithmic design.

Our results show that we can sample a server closer to the target using bounded samples at each node. In order to obtain a node that is β ($\beta \in (0, 1]$) times closer to the target than the current node, we need to uniformly sample enough neighbors from the proximity region of each current node.

Without loss of generality, assume that a node P needs to locate a node Q that is β ($\beta \leq 1$) times closer to a target T , which implies that $d_{QT} \leq \beta \times d_{PT}$. Let $d_{PT} = r$. We can see that node Q must be covered by the ball $B_P(\rho r)$, since $d_{PQ} \leq \rho \max\{d_{PT}, d_{QT}\} = \rho r$. Fig 4 shows an example of sampling a node closer to the target T in the closed ball $B_P(\rho r)$ in the growth dimension.

We first quantify the volume differences of balls with identical centers but different radii.

Lemma VI.1. *Given a ρ -inframetric with growth $\gamma_g \geq 1$, for any $x \geq \rho$, $r > 0$ and any node P , the volume of a ball $B_P(r)$ is at most x^α smaller than that of the ball $B_P(xr)$, where $\log_\rho \gamma_g \leq \alpha \leq 2 \log_\rho \gamma_g$.*

Lemma VI.1 states that the volume differences of the balls with identical centers and different radii are bounded by x^α , where x is the multiplicative ratio between different radii, and the parameter α lies in a bounded interval.

We calculate α by varying the radius r and the multiplicative ratio x as shown in Fig 5. We can see that α is mostly below 1, and decreases close to 0 quickly with increasing radius r or multiplicative ratio x . Therefore, the volume difference x^α scales **sub-linearly** in most cases. On the other hand, for small radius r or low multiplicative ratio x , the volume difference x^α may scale **ultra-linearly**.

Furthermore, we also characterize the inclusion relation of balls with different centers, which generalizes the inclusions of balls around a node pair in the metric space [44]. Lemma VI.2 lays the foundation for uniform sampling nodes to perform DNNS on the inframetric model.

Lemma VI.2. *(Sandwich lemma) For any pair of node p and q , and $d_{pq} \leq r$, then*

$$B_q(r) \subseteq B_p(\rho r) \subseteq B_q(\rho^2 r)$$

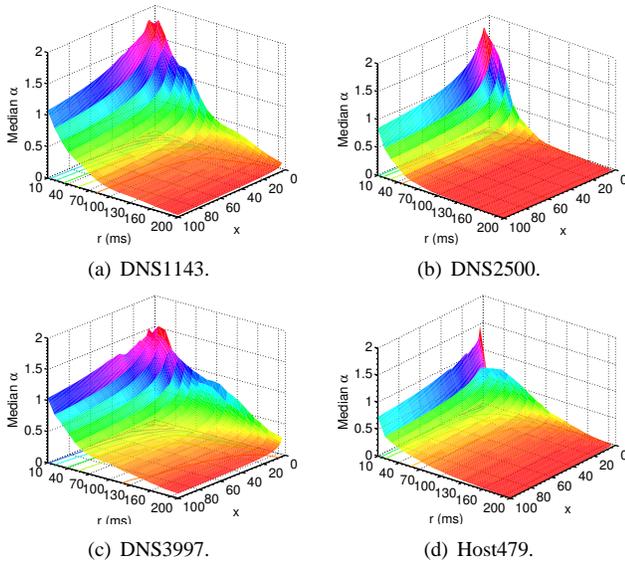


Fig. 5. Median α as function of the radius r and the multiplicative ratio x .

Using Lemma VI.1 and VI.2, we can quantify the size of sampled neighbors, to assure that at least one neighbor lies in the closed ball $B_T(\beta r)$.

Theorem VI.3. (Sampling efficiency in the growth dimension) For a ρ -inframetric model with growth $\gamma_g \geq 1$, for a service node P , and a DNNS target T satisfying $d_{PT} \leq r$, when selecting $3\left(\frac{\rho^2}{\beta}\right)^\alpha$ nodes uniformly at random from $B_P(\rho r)$ with replacement, with probability of at least 95%, one of these nodes will lie in $B_T(\beta r)$, where $\log_\rho \gamma_g \leq \alpha \leq 2\log_\rho \gamma_g$ and $\beta < 1$.

Since α and ρ are determined by the delay space, we can see that the number of samples decreases with increasing delay reduction threshold β . As β approaches 1, the number of required samples becomes approximately $3\left(\frac{\rho^2}{\beta}\right)^\alpha \approx 3\rho^{2\alpha} \in [3\gamma_g^2, 3\gamma_g^4]$ based on Lemma VI.1.

B. DNNS on the Inframetric Model

In this section, we present the analysis of DNNS on the Inframetric model. We will show the search accuracy, search periods and search costs related to a DNNS process. We prove that, by recursively following such sampling conditions, we can locate a server that is $1/\beta$ -approximation to the optimal: the delay from the found server to the target is not bigger than $1/\beta$ times that from the nearest server to the target.

First, we review the goal of each DNNS step using the sampling conditions in Sec VI-A. Assume that a node P wants to locate a node that is β times closer to a target T . The goal of the current DNNS step is to locate a node β times closer to the target than the current node P . To that end, Theorem VI.3 shows that we need to sample up to $3\left(\frac{\rho^2}{\beta}\right)^\alpha$ nodes uniformly at random from $B_P(\rho r)$ with replacement.

Based on the sampling condition in Theorem VI.3, performing DNNS in the growth dimension can be formulated into a

simple DNNS procedure in Definition VI.4.

Definition VI.4 (A simple DNNS method in the inframetric model). *sampling $3\left(\frac{\rho^2}{\beta}\right)^\alpha$ neighbors from the closed ball $B_P(\rho d_{PT})$ at each intermediate node P , forwarding the DNNS request to a next-hop node β times closer to the target than the node P , and stopping at a local minima when we can not find such a next-hop node.*

Furthermore, we can quantify the efficiency of found neighbors based on the above DNNS procedure by Corollary VI.6. As a result, we can locate an approximately optimal nearest neighbor for a target T when β approaches one. Furthermore, the number of required search steps is a logarithm function of the ratio Δ of the maximum delay to the minimum delay in the delay space, indicating that the DNNS queries can complete quickly.

Definition VI.5 (ω -approximation). *For a DNNS request with target T , a found nearest neighbor A is a ω -approximation, if the delay between A to T is smaller than ωd_* , where d_* is the delay between the real nearest neighbor to T .*

Corollary VI.6. *For a relaxed inframetric model with growth γ_g , according to the DNNS process in Definition VI.4, the found nearest neighbor is a $\frac{1}{\beta}$ -approximation, and the number of search steps is smaller than $\log_{\frac{1}{\beta}} \Delta$, where Δ is the ratio of the maximum delay to the minimum delay of all pairwise delays.*

C. Limitations of Theoretical Results

To find a better next-hop neighbor without missing any closer nodes, based on the DNNS analysis in the inframetric model in Sec VI-B, we should sample approximately $3\left(\frac{\rho^2}{\beta}\right)^\alpha$ nodes whose delays to current node P are not larger than ρd_{PT} . However, the number of the candidate neighbors may be quite high, as shown in Fig 6. We can see that the number of required samples exceeds 100 accordingly, for β below 0.4 or α above 1. Such high number of samples implies that we need extremely large number of samples for continuing the DNNS query.

On the other hand, the number of samples decreases with decreasing α or with increasing β . When α is below 1, the number of samples is below 33 if the delay reduction threshold β is above 0.8. As a result, we can see that we need to choose a large β in order to reduce the number of samples, since the median values of α are mostly no more than 1 from Fig 5.

D. Comparison with Previous Inframetric Study

Our relaxed inframetric model is inspired by the seminal study on the inframetric model [43] that assumes the symmetry of the distance function. We extend the inframetric model study for the Internet delays in four aspects:

- We extend the inframetric model to allow both symmetric and asymmetric distance functions, which generalizes the RTTs and OWDs that are important for latency-sensitive applications.

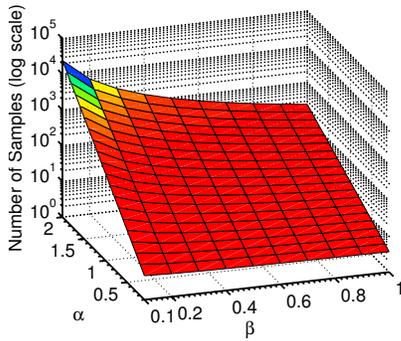


Fig. 6. The number of sampled neighbors $3\left(\frac{\rho^2}{\beta}\right)^\alpha$ by varying the volume difference parameter α from the interval $[0, 2]$ based on the analysis in Sec VI-A and the delay reduction threshold β . We set the inframetric parameter ρ to be 3 to represent most triples.

- We clearly show the relation between inframetric parameter ρ and the TIV. The inframetric parameter $\rho \leq 2$ is a necessary but not sufficient condition for no TIVs.
- We formulate the DNNS problem on the relaxed inframetric model and propose a simple DNNS method that finds approximately nearest neighbor for any target using at most logarithmic search hops. Interestingly, our simple DNNS method works on both symmetric and asymmetric delay metrics.

VII. REALIZING A PRACTICAL DNNS

A. Overcoming Limitations of the Simple DNNS Method

Recall that the measurement costs limits the usefulness of the simple DNNS method defined in Def VI.4 from Sec VII-A. Besides, in the distributed system context, since each service node does not have the global view of the delay space, sampling enough neighbors from the closed ball centered at each service node is difficult. We discuss design principles to tackle these two difficulties in this section.

1) *Reduce Measurement Costs:* We reduce the measurement costs in two complementary approaches: (i) Given that the number of required samples of the simple DNNS method depend on varying parameters, we seek to modify the parameters to obtain the lower bound of the required number of samples. (ii) Given that network coordinates can be used for delay estimations, we avoid complete measurements from selected samples to the target using delay estimations.

First, recall that the number of samples for the simple DNNS method increases quickly with decreasing delay reduction threshold β . Therefore, to reduce the number of samples, we should set the delay reduction threshold β to be close to 1. On the other hand, since the approximation ratio of the simple DNNS method is $1/\beta$, we can see that large β also leads to better approximations of nearest neighbors. As a result, we set β to 1 in order to reduce the number of samples and obtain the best approximation accuracy.

Second, although we reduce the number of samples using modified β , we still need delay measurements between selected samples to the targets, which consume the bandwidth costs and CPU loads of service nodes. Therefore, we hope to reduce the required delay measurements while obtaining

the sample that is closest to the target. To that end, we use *delay estimations based on network coordinates* to reduce the delay measurement costs. However, since the delay estimations incur errors due to the embedding distortions of network coordinates, simply using delay estimations to find the nearest neighbors becomes less reliable. Instead, we issue delay measurements when the delay estimations are inaccurate, so as to avoid the inaccurate delay estimations.

2) *Sample Enough Neighbors For Continuing DNNS Query:* Based on the simple DNNS method, each DNNS service has to maintain enough neighbors covering different delay ranges in the delay space, in order to find the nearest neighbor to any target. Therefore, each node has to maximize its diversity in the neighbor set.

Gossip based neighbor management is frequently used for existing DNNS methods. For example, Meridian [13] and OASIS [9] use an anti-entropy gossip protocol to discover neighbors, and store neighbors using rings of neighbors called concentric rings. However, during our experiments, the innermost and outermost rings in the concentric ring often find no or only few neighbors compared to the capacity of the ring, while the rest of rings with radii lying in the middle portion of the delay distributions are filled with too many neighbors, leading to frequent ring management events, incurring heavy computation and communication overhead.

We explain the insufficiency of the gossip process in details. Assuming that we know the complete delay matrix, for each node, we compute the percent of mapped nodes for each ring, which serves as an upper bound of sampled neighbors for that ring. Then we can analyze whether the distributions of mapped nodes in concentric rings affect the gossip process. As shown in Fig 7, we can see that most nodes are mapped into a few number of rings, whose delay ranges lie in the middle portion of the delay distributions. However, only quite a few nodes are mapped into the innermost and outermost rings, which result in a skewed distribution of mapped nodes for the concentric rings. As a result, since the gossip process adopts the uniform sampling approach, the gossip process will inevitably sample insufficient neighbors from those rings that have too few mapped nodes.

Accordingly, to improve the concentric ring maintenance, we need to sample enough neighbors that lie in different delay ranges. To that end, we propose to find nearest neighbors and farthest neighbors for each service node, in order to fill the innermost and outermost rings in the concentric ring.

B. Our Design

Based on the design principles in Sec VII-A, we design a novel DNNS method named *HybridNN* (Hybrid Nearest Neighbor Search). We present an overview of HybridNN. To sample enough candidate neighbors from the proximity region of the current node, each node must first maintain a neighbor set that contains enough neighbors within each proximity region. Then using the neighbor set, we select candidate neighbors using the sampling conditions of the simple DNNS method, in order to cover the neighbors closer to the target with high probability. Next, we determine the candidate

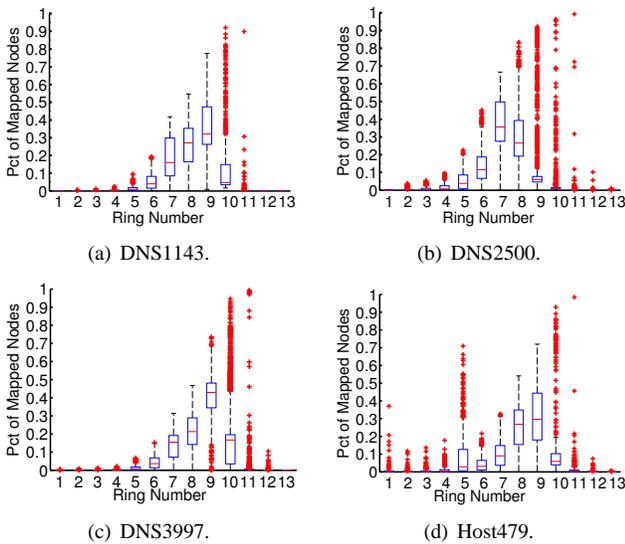


Fig. 7. The percent of mapped nodes into different rings, assuming that we obtain the complete delay matrix. The i -th ring contains neighbors whose delays to a node P lie in the interval $(\alpha s^{i-1}, \alpha s^i]$, with $i > 0$, α a constant, s a multiplicative increase factor ($\alpha = 1$, $s = 2$ ms as configured by Wong et al. [13]). Besides, since our objective is to determine the distribution of nodes mapped into the concentric ring, we do not limit the maximum capacity of each ring.

neighbor closest to the target, using delay estimations and direct probes, in order to obtain a better tradeoff between sampling bandwidth and accuracy. Finally, using the currently nearest candidate neighbor to the target, we determine whether to terminate the DNNS query. As shown in Fig 8, HybridNN is composed of five components:

Neighbor Maintenance: This component maintains the neighbor set for DNNS queries. Since nodes are mapped into the rings at the middle portion of the concentric ring, which implies that neighbors mapped into the head portion and tail portion of the concentric ring are difficult to be sampled using the uniform sampling based approach. As a result, we need to increase the sampling probability of such neighbors, in order to fulfill the sampling conditions for DNNS queries. To that end, we over-sampling neighbors in the head portions and tail portions of the concentric rings, besides we uniformly sampling neighbors located in the middle portions of delays and.

Selecting Candidate Neighbor: This component selects candidate neighbors to satisfy the sampling conditions of the simple DNNS method. When a node P receives a DNNS query, node P determines its delay towards the target T , then selects neighbors from its diversity-optimized neighbor sets (Sec VII-C) by covering possible closer neighbors towards the target T (Sec VII-D). Furthermore, we prune those neighbors that could mislead the DNNS query into poor local minima.

Coordinate Maintenance: This component updates the coordinate of the target in order to estimate delays to targets from candidate neighbors, since the target machine may not have the coordinate for delay estimation (Sec VII-E). Additionally, each service machine maintains a network coordinate used for delay estimations.

Determining Closest Neighbor: This component determines

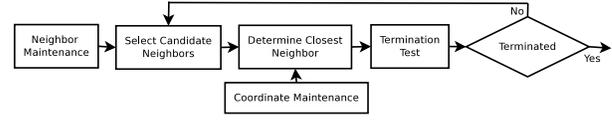


Fig. 8. The flow chart of four search steps at a service node for a DNNS query.

the neighbor nearest to the target (Sec VII-F). Each node computes the candidate neighbor closest to the target using delay estimations and direct probes, in order to balance between the measurement costs and measurement accuracy.

Termination Test: This component determines to continue or stop a DNNS query (Sec VII-G). Recall that in previous section we set the delay reduction threshold β to be 1 in order to reduce the number of samples and obtain better approximation ratios to the optimal results. Therefore, HybridNN conservatively terminate the DNNS query only when all candidate neighbors having larger delays than the current node.

Finally, HybridNN uses an extensible delay measurement interface. For instance, by default HybridNN simply use the system-built-in Ping command to obtain pairwise RTT measurements. When there exist an on-demand OWD probe service such as Reverse Traceroute [18], HybridNN configures a RPC interface to request the pairwise OWD results.

C. Neighbor Maintenance

In order to facilitate the neighbor sampling for DNNS forwarding, each service node maintains neighbors that are sampled from different regions in the delay space. We introduce the neighbor discovery and update in this section.

1) *Organize Neighbors Into Rings for Proximity Selection:* Since the proximity region for neighbor sampling in the simple DNNS method is a closed ball, we choose the concentric ring to organize neighbors for each node. For instance, if we need to locate all neighbors that are at most d_2 ms away, we select all neighbors from those rings whose ring numbers are at most $\lceil \log_2 d_2 \rceil$.

An important parameter for the concentric ring is its ring size Δ , which determines the maximum number of neighbors per ring. Since we need to sample enough neighbors using the concentric ring to guarantee to locate a neighbor closer to the target with a high probability, we analytically determine the choice of Δ as follows. First, the total number of samples $3\left(\frac{\rho^2}{\beta}\right)^\alpha$ is within the interval $[3\gamma_g^2, 3\gamma_g^4]$, since we set the delay reduction threshold β to 1. Therefore, if we set the number of neighbors Δ at each ring to be at least $O(\gamma_g^2)$, we can ensure that with a high probability, we can find a neighbor that is closer to the target than the current node P . Furthermore, since γ_g is low on average from previous sections, we can set the number of neighbors Δ to be a modest integer (8 by default).

Furthermore, to adapt to the dynamics of delays, we use a moving median as a latency filter for extracting stable delay measurements to each neighbor [45], which allows to have up-to-date delay estimates resilient to the measurement noises.

2) *Biased Sampling based Neighbor Discovery*: Based on the distribution of neighbors for each ring in the previous section, we have seen that we need to over-sample neighbors mapped into the head portion and the tail portion of the concentric rings. To that end, we adopt both uniform sampling and over-sampling approaches.

Uniform sampling. We reuse the gossip process in Meridian. Briefly, each node P periodically starts the gossip process by uniformly selecting a neighbor Q from P 's concentric ring as communication partner, and sends a gossip request message to node Q containing randomly sampled neighbors, one neighbor per non-empty ring. When Q receives the gossip request, Q will send a gossip ACK to P immediately; besides, Q iteratively sends gossip requests towards the sampled neighbors in the gossip request message of P .

Finally, if we use the RTT metric, then node P inserts Q into the corresponding ring according to the round trip delays measured as the period between the gossip request and the gossip ACK. Alternatively, if node P is able to measure the one-way delay from P to Q , then node Q is inserted into the corresponding ring according to the one-way delay from P to Q .

Over-sampling. Our goal is to sample enough neighbors from those mapped neighbors lying in the head and tail portions of the concentric rings. For this purpose, we use K closest neighbor search and K farthest neighbor search. The returned nodes are directly stored into the concentric ring, as the delay values between the current service node to the returned nodes are obtained during the K closest neighbor search and K farthest neighbor search processes.

- *K closest neighbor search.* Each node P periodically finds nearby nodes by issuing K closest neighbor search with itself as target. Here K is a system parameter. Firstly, node P randomly selects a neighbor Q from its concentric ring, and sends to Q a K nearby neighbor search message. Then node Q starts a K closest neighbor search process. After the K closest neighbor search process is completed, found nearby nodes and the corresponding delays to P are returned to node P , and P saves these returned nearby nodes into its concentric ring.
- *K farthest neighbor search.* Similar as the K closest neighbor search process, each node P periodically issues K farthest neighbor search. Later, the K farthest neighbor search results include found distant neighbors and the corresponding delay values to node P . P stores the returned distant neighbors into its concentric ring by the corresponding delay values.

Due to space limits, the details for K closest neighbor search and K farthest neighbor search are omitted here, which can be found in the full technical report [36].

3) *Replacing Suboptimal Neighbors Without Probes*: In order to bound the memory overhead of the concentric ring, we need to manage the size of the concentric rings when some rings reach their maximum capacity Δ . To reduce CPU costs due to frequent ring managements, we lower the frequency of ring managements: we first set up another tolerance threshold Δ_t for each ring; then we begin the ring management when some rings having at least $\Delta + \Delta_t$ neighbors; during the ring

management, we remove Δ_t neighbors from those rings that have at least $\Delta + \Delta_t$ neighbors.

When we need to remove Δ_t neighbors from some rings, we follow the removing philosophy of Meridian: preserve those that maximize the diversity of neighbors in a ring using the maximal hypervolume polytope algorithm ([13]). This is because the higher diversity in the neighbor set translates to better chances of locating a nearby nodes for any target. However, the maximal hypervolume polytope algorithm requires all-pair delay measurements of nodes in a ring, which needs $O(\Delta^2)$ probes. In order to avoid such measurements, we turn to adopt network coordinates for delay predictions.

For delay predictions, we use the revised Vivaldi algorithm [21] that is robust to TIVs [20]. We denote the revised Vivaldi [20] as *TIV-Vivaldi*($x_i, e_i, d_{ij}, x_j, e_j$), where the input x_i, x_j denote the coordinate of node i and j , respectively; the input e_i, e_j denote the averaged error of node i 's and j 's coordinates, respectively. The output of *TIV-Vivaldi* are the updated coordinate x_i and coordinate error e_i of node i .

Each service node passively maintains a coordinate, and estimates delays using coordinate distances. Besides, for estimating delays with neighbors in the concentric ring, each service node also stores its neighbors' coordinates.

Since delay varies, each node updates its own and cached coordinates periodically. Rather than introduce additional delay probes, we update coordinates by reusing the delay measurements to other service nodes during the biased sampling procedure. Therefore, we significantly reduce the maintenance costs compared to Meridian. First, each node receiving the gossip message piggybacks its coordinate to the sender along with the acknowledged gossip message. After receiving the coordinate from the gossip receiver node, the gossip sender node stores the new coordinate of the gossip receiver node, and updates its own coordinate by triggering *TIV-Vivaldi* using the delays obtained during the gossiping process.

D. Select Candidate Neighbors

Assume that node P receives a DNNS query to the target T . Based on the sampling conditions of the simple DNNS method, node P needs to select $3\left(\frac{\rho^2}{\beta}\right)^\alpha$ neighbors whose delays to node P are in the delay range $[0, \rho d_{PT}]$. Since each ring contains $O(\gamma_g^2)$ neighbors, we simply select all neighbors of rings numbered in the range $[1, \lceil \log_2(\rho d_{PT}) \rceil]$ as candidate neighbors.

Furthermore, we also prune several neighbors that mislead the DNNS process. First, candidate neighbors that contain too few non-empty rings are more likely to provide no hints on continuing the DNNS queries, thus the DNNS queries can be trapped into local minima, due to the neighbors' sparse diversity of the delay space. Therefore, we remove all neighbors with fewer than τ non-empty rings ($\tau = 4$ by default). Second, all neighbors that have received the identical DNNS query should be removed in order to avoid the search loops. Therefore, let the *forwarding path* of a DNNS query be the sequence of nodes forwarding the query. we remove any node on the forwarding path.

E. Coordinate Maintenance for Targets

In order to reduce the delay measurement costs, we predict delays from service nodes to the target, since each service node has computed its network coordinate during the neighborhood management process (Sec VII-C3). As a result, reusing the coordinates for predicting delays can reduce the measurement costs.

Unfortunately, we may not know the coordinate of the target, as the target can be any machine on the Internet. Therefore, we propose to compute the coordinate for the target on-the-fly based on the *TIV-Vivaldi*.

First, when node P receives the DNNS query for a target T , node P will initialize the network coordinate x_T for target T if T 's coordinate is not stored in the DNNS query message. To that end, node P asks a fixed number of neighbors (at most 10) to directly probe the target T . Then, node P updates target T 's coordinate by *TIV-Vivaldi* using the coordinates and delay measurements from these neighbors to target T , which updates T 's coordinate x_T and coordinate error e_T as the output of *TIV-Vivaldi*. Finally, node P stores target T 's coordinate into the DNNS query and forwards to the next-hop node for recursive search. This completes the coordinate initialization for the target T .

Second, after initializing T 's coordinate, each node Q that forwards the DNNS query will update target T 's coordinate for better convergence of target T 's coordinate. To that end, each node Q applies *TIV-Vivaldi* to update target T 's coordinate x_T and coordinate error e_T , using node Q 's coordinate and delay d_{QT} the target T .

F. Determine Closest Neighbor

After we assign a network coordinate to the target in Sec VII-E, we can use the network coordinate distances to approximate the real-world delay and reduce the measurement costs. Nevertheless, since the coordinate distances are only approximations, closest neighbors selected according to the network coordinates may be inconsistent with the real ones.

Therefore, we locate closest neighbors to the target T from the candidate neighbors found in Sec VII-D, by combining the delay predictions with a small number of direct probes.

First, based on the coordinate distances from candidate neighbors to target T , we find top- m nearest neighbors S_c to the target T from the candidate neighbors.

Second, since coordinate distances may be erroneous, we also choose those candidate neighbors S_e whose coordinates are not reliable. Since each *TIV-Vivaldi* coordinate x_i is accompanied by a coordinate error metric e_i [20], we choose unreliable neighbors whose coordinate errors exceed a threshold. We found that setting the threshold to be 0.7 can significantly reduce the negative impact due to the coordinate inaccuracies.

Third, to adapt to coordinate errors caused by TIV, since high coordinate distance errors indicate violations of triangle inequality [20], we simply include all candidate neighbors S_t whose coordinate distance and real delay towards the current node P differs by more than 50 ms, which has good tradeoff between accuracy and bandwidth costs.

Finally, using the union of selected candidate neighbors $S_* = S_c \cup S_e \cup S_t$, the current node P asks neighbors in S_* to probe the delays to target T , from which node P determines the closest neighbor. Ties are broken by choosing the neighbor with most accurate coordinate.

G. Termination Test

Recall from Sec VII-A, HybridNN set the delay reduction threshold β to be 1, in order to reduce the number of selected neighbors and obtain better approximation ratios for the found nearest neighbors. Therefore, when the closest neighbor selected from Sec VII-F has a larger delay to the target than that of the current node P , node P terminates the DNNS query. Then node P sends the currently closest node to the host that issues the DNNS query.

VIII. EXTENSIONS TO HYBRIDNN

HybridNN can be readily extended to search more than just one nearest node. Here we will just give two examples namely, K closest neighbor search and K farthest neighbor search, which are both utilized to oversample neighbors in the network delay space in order to increase the diversity for neighborhood management.

A. K Distributed Nearest Neighbor Search

The K Distributed Nearest Neighbor Search (KDN²S) aims to locate the K nearest neighbors to a target T , where K is a system parameter. To store the found nearest neighbors, we append a new field $M.\Omega$ that caches nearest neighbors to the DNNS query message M .

A naive KDN²S solution is based on the *finding and removing* approach: first we find one closest neighbor towards the target based on the HybridNN algorithm, then we delete the found nearest neighbor from the system, and we restart the HybridNN algorithm from the same query node until we locate K nearest servers to the target. Nevertheless, deleting the closest neighbors from the system is not practical for a large-scale system due to the broadcasting communication overhead, and repeated DNNS processes increase the query overhead for the service nodes on the DNNS forwarding paths.

On the other hand, if we assume that the concentric ring of each node does not append new neighbors, the network coordinate of each node keeps unchanged and the network delays keep stable during the period of a KDN²S query, we find that there exists *temporal correlation* in the forwarding paths of consecutive DNNS queries starting from the identical node in the naive KDN²S solution: *if we issue a new DNNS query from the same starting node immediately after the preceding DNNS query, then the forwarding path truncated the last-hop node of the new DNNS process is a subpath of the forwarding path of the preceding DNNS query, since we can see that the intermediate nodes on these two forwarding paths are identical in HybridNN*. Our assumption generally holds after the network coordinates converge and the concentric rings contain enough neighbors. Furthermore, the constancy of end to end network delays has been confirmed to be on the orders

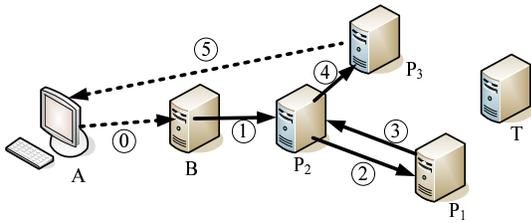


Fig. 9. KDN²S.

of hours by Zhang and Duffield [46] as well as the iPlane project [33], [34].

Using the temporal correlation of consecutive forwarding paths from the same starting node, we propose a backtracking based KDN²S algorithm, as shown in Algorithm 1. After we find one nearest neighbor and terminate at a service node P_1 by HybridNN, we resume the KDN²S query from P_1 , by backtracking from P_1 to its predecessor node P_2 on the DNNS forwarding path, and by recursively finding the nearest neighbor at P_2 , until we locate K nearest neighbors. With backtracking, the KDN²S resumes the query at service nodes that are close to the target, therefore we can quickly locate new nearest neighbors with reduced forwarding overhead compared to the naive KDN²S solution.

Fig 9 gives an example of KDN²S using Algorithm 1. Suppose an end host A needs two nearest neighbors to the target T . Node A sends a KDN²S request to a service node B . Then B starts the KDN²S by forwarding a KDN²S query M to a neighbor P_2 closer to T . Similarly, P_2 forwards the query M to P_1 . Now node P_1 finds that it cannot find a neighbor closer to the target T than itself, therefore, P_1 is the first nearest neighbor to the target. Then P_1 appends its address into $M.\Omega$ as a found nearest neighbor. Next P_1 triggers the KDN²S backtracking step by forwarding M to P_1 's predecessor P_2 on the KDN²S forwarding path. On receiving M , P_2 excludes P_1 from the choice of candidate neighbors, and finds a new neighbor P_3 closer to the target T than P_2 . Then P_2 forwards M to P_3 . P_3 decides that it is the closest node to T among its neighbors. Therefore, P_3 appends itself to $M.\Omega$ as a new nearest neighbor. Finally, P_3 sends the found nearest neighbors in $M.\Omega$, i.e., P_1 and P_3 , to the end host A , which completes the KDN²S.

B. K Distributed Farthest Neighbor Search

Similar as the KDN²S, K Distributed Farthest Neighbor Search (KDFNS) is also based on the backtracking idea. First, we locate one farthest neighbor and terminate at a service node P , then we backtrack from node P to its predecessor node on the forwarding path to recursively locate the rest $K - 1$ farthest neighbors.

To locate one farthest neighbor, we recursively forward the KDN²S query to a service node P_1 that is at least $(1 + \beta_{farthest})$ ($\beta_{farthest}$ is 1.2 by default) times farther to the target T than the current service node P . In other words, we need to locate a node that is not covered by the ball $B_T((1 + \beta_{farthest})d_{PT})$. Since $B_T((1 + \beta_{farthest})d_{PT}) \subseteq B_P(\rho(1 + \beta_{farthest})d_{PT})$ by

Algorithm 1: The pseudo-code of KDN²S.

- 1: $KDN^2S(H, T, K, M)$
 - 2: {Input: current node H , the target T , required number of closest neighbors K , query message M }
 - 3: {Output: nearest neighbors to T }
 - 4: **if** $|M.\Omega| == K$ **then**
 - 5: Return $M.\Omega$; {enough closest neighbors}
 - 6: **end if**
 - 7: $S \leftarrow \text{chooseCandidates}(P, T, M)$;
 - 8: $S \leftarrow S - M.\Omega$ {remove found nearest neighbors to avoid search loops}
 - 9: $x_T \leftarrow \text{InitTargetCoord}(P, T)$;
 - 10: $[u_1, S_c, D_T] \leftarrow \text{NearestDetector}(P, S, x_T, M)$;
 - 11: $[\phi_1, d_{\phi_1 T}, P_1] \leftarrow \text{TerminateTest}(P, u_1, S_c, D_T, M)$; {find one closest neighbor, and terminate at node P_1 }
 - 12: $M.\Omega \leftarrow M.\Omega \cup \{\phi_1\}$; {cache ϕ_1 into the query message}
 - 13: Select the predecessor node P_2 of node P_1 on the forwarding path $M.\text{Path}$; {find the predecessor for backtracking}
 - 14: $KDN^2S(P_2, T, K, M)$; {recursive search}
-

the sandwich lemma in Lemma VI.2, P_1 needs to be at least $\rho(1 + \beta_{farthest})d_{PT}$ from node P .

Accordingly, in each search step, we try to find such node P_1 from the concentric ring of the current service node P , whose delay value to P is larger or equal the $\rho(1 + \beta_{farthest})d_{PT}$. If there exists a such node P_1 , then node P_1 recursively runs the KDFNS as node P . Otherwise, if we can not locate such node P_1 , the search is terminated, and the currently farthest node to the target is cached as a farthest neighbor to the target. Afterwards, we select the rest $K - 1$ distant neighbors by the backtracking process similar as that in K closest neighbor search.

Algorithm 2 shows the complete KDFNS process. First, we choose candidate neighbors satisfying the delay constraint to the current service node P . Then we find the farthest neighbor to the target (*FarthestDetector*()) combining the delay predictions with direct probes in order to reduce the measurement overhead. Specifically, we choose m farthest neighbors from the candidate neighbors; besides, we also add neighbors with uncertain coordinates and erroneous predictions similar as Sec VII-F. Next, we determine one farthest neighbor recursively (*FarthestTerminateTest*). Finally, from the terminating node P_1 , we backtrack to the predecessor node of P_1 on the forwarding path, and recursively run the KDFNS until we locate enough farthest nodes to the target.

IX. SIMULATION

In this section, we report the results of simulation experiments based on the real-world data sets in Sec IV.

A. Experimental Setup

We compare HybridNN with several DNNS algorithms. (1) **Vivaldi**. We compute the coordinate of each node based on the Vivaldi algorithm [45], and find the nearest service nodes for each requesting node using shortest coordinate distances. The coordinate dimension for Vivaldi is 5. (2) **CoordNN**. To quantify the usefulness of direct probes of HybridNN, we present a DNNS algorithm CoordNN, which is identical with

Algorithm 2: The pseudo-code of KDFNS.

```

1:  $KDFNS(H, T, K, M)$ 
2: {Input: current node  $H$ , the target  $T$ , required number of
   farthest neighbors  $K$ , query message  $M$ }
3: {Output: farthest neighbors to  $T$ }
4: if  $|M.\Omega| == K$  then
5:   Return  $M.\Omega$ ; {complete the KDFNS}
6: end if
7:  $S \leftarrow \text{chooseFarthestCandidates}(P, T, M)$ ; {choose neighbors
   whose delay values to  $P$  is larger than or equal to
    $\rho(1 + \beta_{farthest})d_{PT}$ }
8:  $S \leftarrow S - M.\Omega$  {remove found farthest neighbors to avoid
   search loops}
9:  $x_T \leftarrow \text{InitTargetCoord}(P, T)$ ;
10:  $[u_1, S_c, D_T] \leftarrow \text{FarthestDetector}(P, S, x_T, M)$ ; {select the
   farthest neighbor to  $T$  from  $S$ }
11:  $[\phi_1, d_{\phi_1 T}, P_1] \leftarrow \text{FarthestTerminateTest}(P, u_1, S_c, D_T, M)$ ;
   {find one farthest neighbor, and terminate at node  $P_1$ }
12:  $M.\Omega \leftarrow M.\Omega \cup \{\phi_1\}$ ; {cache  $\phi_1$  into the query message}
13: Select the predecessor node  $P_2$  of node  $P_1$  on the forwarding
   path  $M.\text{Path}$ ; {find the predecessor for backtracking}
14:  $KDFNS(P_2, T, K, M)$ ; {recursive search}

```

TABLE I
PARAMETER VALUES OF HYBRIDNN FOR SIMULATION.

Parameter	Meaning	Value
Δ	maximal size of the ring	8
$\Delta + \Delta_t$	threshold of the ring size for ring updates	10
β	nearest search threshold	1
ρ	inframetric parameter	3
$ x $	coordinate dimension	5
K	size of sampled neighbors for neighbor discovery	10
m	number of neighbors for direct probes	4
τ	number of non-empty rings	4

HybridNN except that it uses only and no direct probes when determining the best next-hop neighbors. (3) **DirectDN2S**. To evaluate HybridNN, we present a DNNS algorithm DirectDN2S, which is identical with HybridNN except that it only utilizes direct probes for finding next-hop best neighbors without pruning neighbors based on coordinate distances as HybridNN. (4) **Meridian** [13]. Meridian recursively forwards the DNNS queries to a node that is β times closer to the target than the current node, and returns the found nearest neighbor when no such node is selected. We configure the parameters of Meridian algorithm identical with the original configuration by Wong et al. [13], with the delay reduction threshold β as 0.5, the upper bound on the size of each ring as 10, and the number of rings in the concentric ring is 20.

For HybridNN, the default configuration is summarized in Table I. CoordNN and DirectDN2S share identical parameters with HybridNN. We also evaluated the sensitivity of parameters for HybridNN, which is reasonably robust against the parameter choices. The detailed sensitivity results of system parameters for HybridNN can be found in the technique report published online [36].

We have developed a discrete-time simulator for DNNS. The simulator randomly chooses a set of nodes as service nodes (by default 500) that can receive DNNS queries. Other nodes in the system are clients that can issue DNNS queries to these service nodes. For Host479, 200 nodes are the service

nodes. The DNNS queries are repeated 10,000 times. For each DNNS query, we uniformly select one client as the target machine, and a random service node receiving the query. Besides, the simulation is repeated 5 times by shuffling the set of service nodes to avoid biases in choosing service nodes. For HybridNN, CoordNN, DirectDN2S and Meridian, the inter-gossip events for neighborhood discovery are generated by an exponential distribution with expected value of 1 second. The inter-ring management events are generated by an exponential distribution with expected value of 2 seconds. For HybridNN, DirectDN2S and CoordNN, the time interval between two oversampling events of K closest neighbor search and K farthest neighbor search are generated by an exponential distribution with expected value of 60 seconds. The inter-DNNS event generation follows an exponential distribution with expected value of 60 seconds. For Vivaldi, the coordinate of each node is updated for 1000 rounds, by uniformly selecting a service node as the counterpart during each round.

The performance metrics for each DNNS query include: (1) **Absolute Error**: defined as the absolute difference between the estimated nearest neighbor j and the real nearest neighbor i to the target T , i.e., $d_{jT} - d_{iT}$. (2) **Relative Error**: defined as the ratio of the absolute error for the estimated nearest neighbor j to the delay between the real nearest neighbor i and the target T , i.e., $\frac{d_{jT} - d_{iT}}{d_{iT}}$. The absolute error quantifies the increased delay values of the estimated nearest neighbors, while the relative error measures the multiplicative ratios to the optimal delay values for the estimated neighbors. Therefore, large relative errors do not necessarily correspond to high absolute errors. (3) **Search Hop**: defined as the number of service nodes on the forwarding path minus one. Therefore, if node A forwards a DNNS query to node B and node B returns the nearest neighbor to the query host, the search hop for the DNNS query is one.

B. Comparison

Absolute Error. Fig 10 shows the absolute errors of the different algorithms. DirectDN2S achieves lowest absolute errors except for the Host479 data sets. HybridNN is close to DirectDN2S in terms of reducing absolute errors, however, HybridNN is the most accurate on Host479 data sets. Next, CoordNN is worse than both DirectDN2S and HybridNN. The accuracy of DirectDN2S and HybridNN compared to CoordNN indicates that utilizing direct probes greatly reduces the inaccuracy of the estimation, while using coordinate distances alone can lead to a bad local minima.

The inaccuracy of DirectDN2S compared to HybridNN on the Host479 data set is rather counter-intuitive. The inaccuracy of DirectDN2S may be caused by the asymmetry in the delay data sets that misleads the greedy search into a local minima, since DirectDN2S is more accurate than HybridNN on the other three data sets that are all symmetric for pairwise delays. On the other hand, HybridNN does not always choose the neighbor closest to the target as the forwarding node, since HybridNN also incorporates the approximated delay predictions when choosing neighbors, which can help HybridNN bypass the bad local minimum caused by the asymmetry in the delay values.

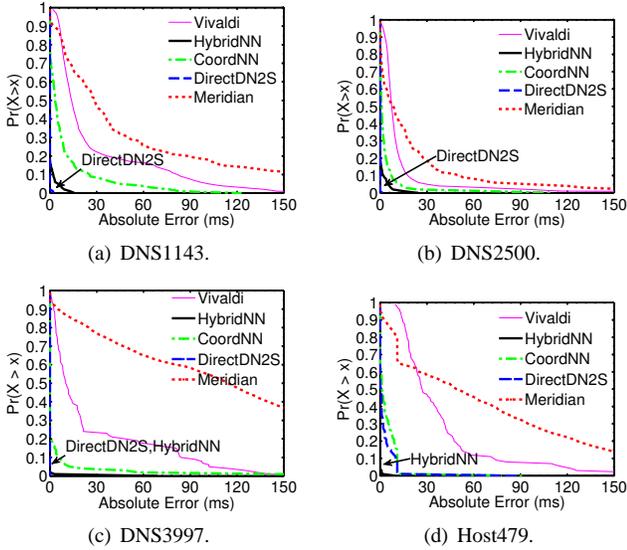


Fig. 10. The CCDFs of absolute errors.

Furthermore, Meridian shows greater absolute errors compared to other algorithms including Vivaldi, which implies that the coordinate distances are at least effective if used in the centralized approach. We are aware that the superiority of Vivaldi over Meridian in most cases are consistent with the experiments independently performed by Choffnes and Bustamante [42]. The main reasons for the less accuracy of Meridian are the local minima caused by the TIV and clustering in the delay space. On the other hand, Vivaldi can adapt to TIV using adaptive coordinate movements.

Relative Error. Fig 11 shows the relative errors of DNNs algorithms. The results are consistent with those of the absolute errors. DirectDN2S achieves near-zero relative errors for most DNNs queries on all data sets except Host479. HybridNN and DirectDN2S have similar accuracy, while HybridNN is more accurate than DirectDN2S on Host479. Furthermore, CoordNN is less accurate than HybridNN, while Meridian and Vivaldi are less accurate than DirectDN2S, HybridNN and CoordNN.

Search hops. Next, we quantify the distributions of the number of search hops for DNNs algorithms, as shown in Fig 12. Recall that the search hops are equal to the lengths of DNNs forwarding paths minus one.

We can see that the search hops of most DNNs queries are rather modest for all DNNs algorithms. Meridian in about 80% of the cases has 2 search hops. While HybridNN and DirectDN2S in over 80% of the cases have no more than 3.

Moreover, almost all searches for Meridian, HybridNN, DirectDN2S are below 6 search hops. On the other hand, CoordNN has longer search hops than Meridian, HybridNN and the DirectDN2S; and a fraction of search hops even exceed 10 on all data sets.

C. Sensitivity of Parameters

In this section, we evaluate the robustness of HybridNN to the system size as well as the choices of system parameters.

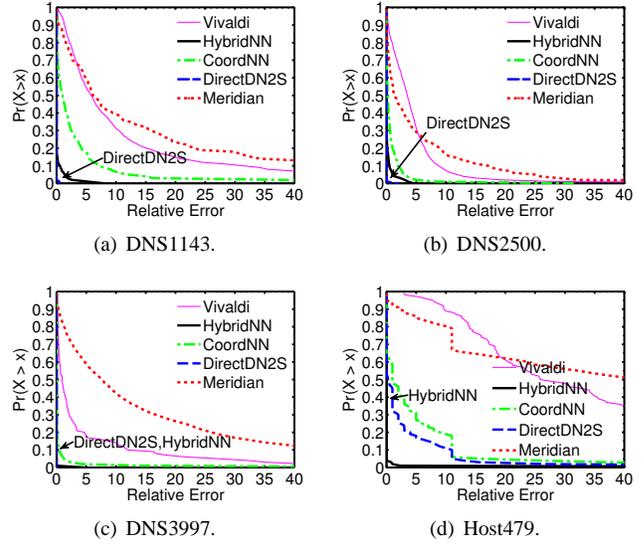


Fig. 11. The CCDFs of relative errors.

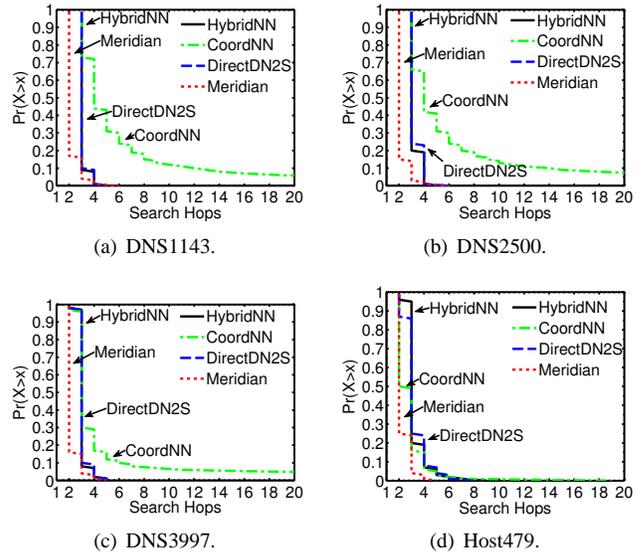


Fig. 12. The CCDFs of search hops.

1) *System Size N*: To evaluate the size of service machines on the performance of HybridNN, we evaluate the performance of HybridNN by increasing the size of service machines. We select target machines randomly from all nodes, including the clients and the service machines, as the size of clients shrinks when increasing the percentage of service machines. Fig. 13 shows the performance of HybridNN with increasing the percentage of service nodes. HybridNN achieves similar accuracy when the size of service nodes increase compared to clients. Therefore, HybridNN is quite robust to the different scales of systems. On the other hand, the query loads of HybridNN increase slowly, for example, HybridNN nearly double the loads when the percentage of service nodes reaches 1.

2) *Inframetric ρ*: Fig. 14 shows the accuracy and loads as the increment of Inframetric parameter ρ. The accuracy of HybridNN is insensitive to choices of ρ. This is because for

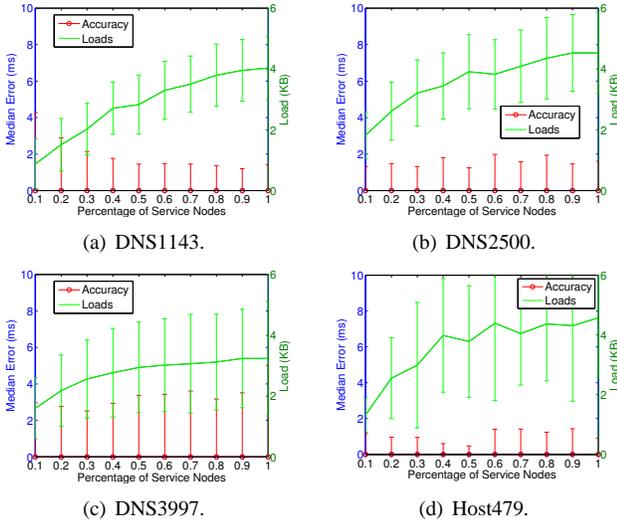


Fig. 13. Size of Service Nodes.

most delays, its ρ -edge metrics are quite lower. Therefore, with lower ρ we can cover possible best next-hop neighbors for DNNs queries. Furthermore, although larger ρ increases the size of possible next-hop candidate neighbors, the loads of DNNs queries of HybridNN keep stable for different ρ , due to that we use nearly constant-sized next-hop nodes. Besides, we can see the standard deviations of errors are quite low for most data sets.

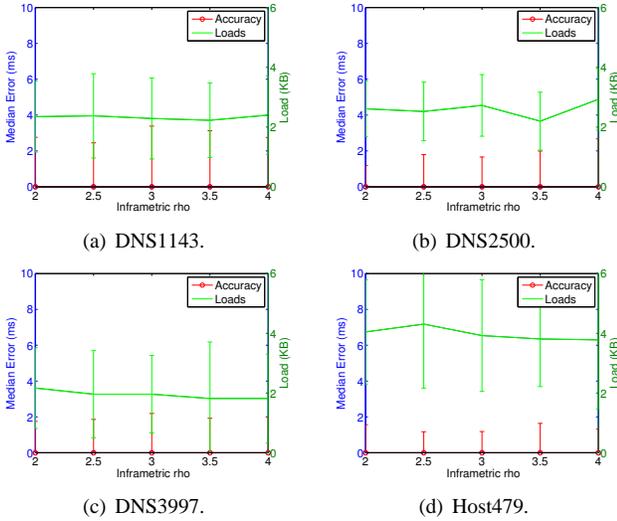


Fig. 14. Inframetric ρ .

3) *Non-Empty Threshold* τ : Fig. 15 shows the accuracy and loads as the increment of Non-empty thresholds for pruning candidate neighbors for next-hop nodes. As the increment of non-empty thresholds for pruning candidate neighbors that have too few rings containing nodes, the standard deviation of HybridNN is reduced before the threshold reaches 4, then increases after the threshold is over 4, and the median errors are increased when the non-empty threshold exceed 8. Besides, the loads are reduced when the non-empty thresholds increase. Therefore, selecting modest-sized non-empty thresholds (e.g.,

4) can keep accuracy and reduce loads.

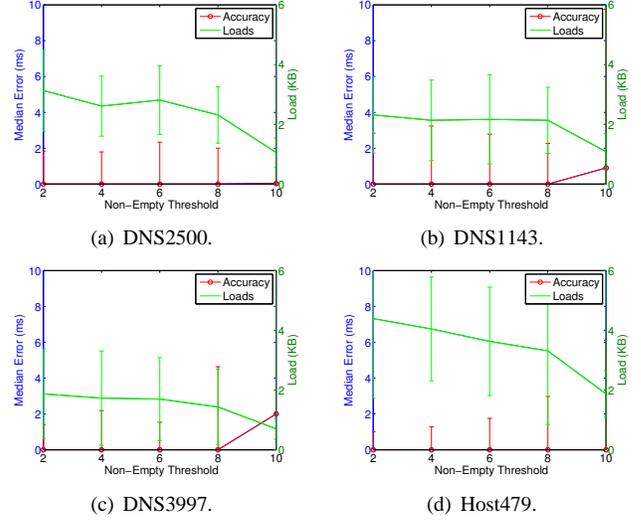


Fig. 15. Non-Empty Threshold.

4) *Coordinate Dimension* $|x|$: Fig. 16 illustrates the accuracy and loads when the coordinate dimension changes. HybridNN achieves similar accuracy and loads as the accuracy of coordinates keeps stably accurate as the dimension is over 3. Therefore, HybridNN can adapt to inaccuracy of different dimensions of coordinates without increasing DNNs query loads efficiently.

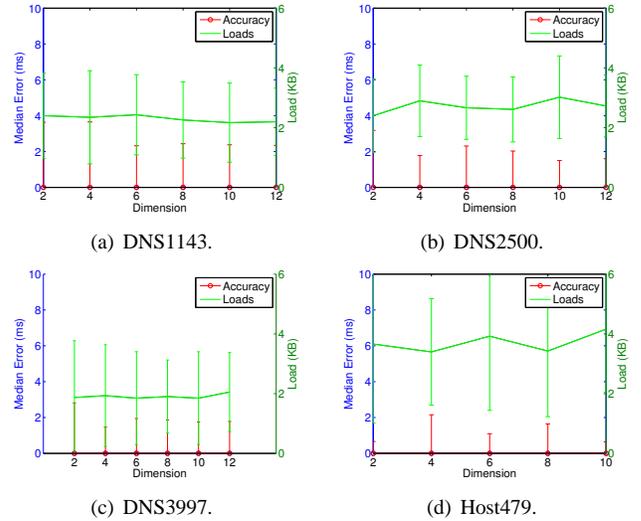


Fig. 16. Coordinate Dimension.

5) *Nodes Per Ring* Δ : Fig. 17 describes the performance of HybridNN with increasing upper bounds of nodes per ring. HybridNN achieves high accuracy event the size of one ring is as small as 5. This is because HybridNN selects neighbors from broader range $[0, \rho d]$, where d is the delay from current node to targets. Besides, the loads of HybridNN grow slowly as the size of ring increases. As HybridNN utilizes coordinate distances to select limited number of candidate neighbor.

6) *OverSampled nearest and farthest nodes* K : Fig. 18 illustrates the performance of HybridNN as the variation

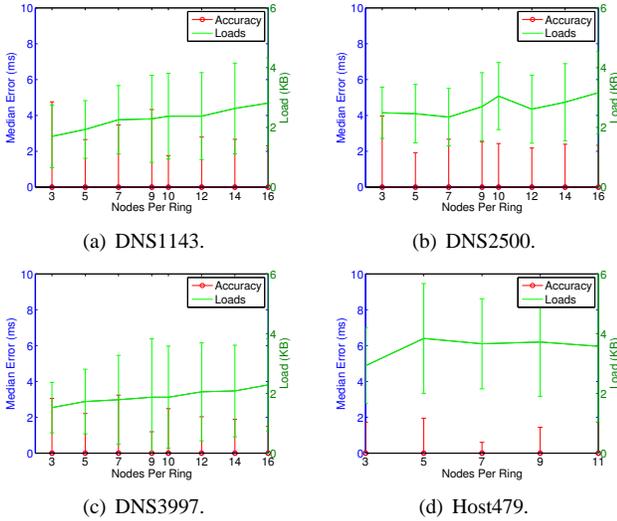


Fig. 17. Nodes Per Ring.

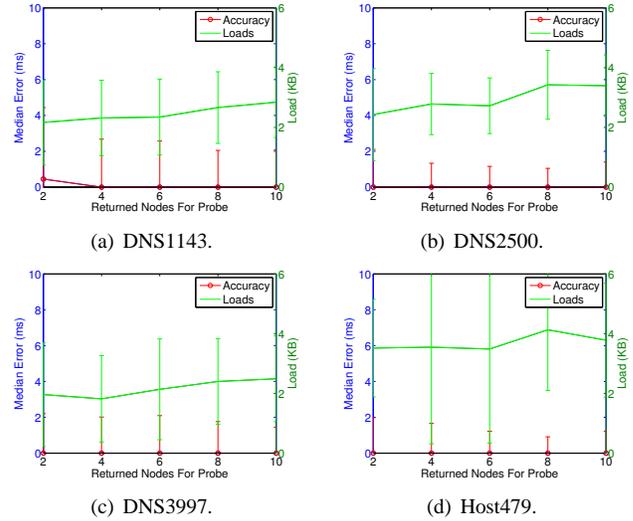


Fig. 19. Returned Nodes For Next-Hop Probes.

of oversampled number of nearest and farthest nodes K . HybridNN achieves similar accuracy and loads when the oversampled size K of nearest neighbors and farthest neighbors. This is because we periodically start the oversampled process, which can find many nearby or far-away nodes accumulatively.

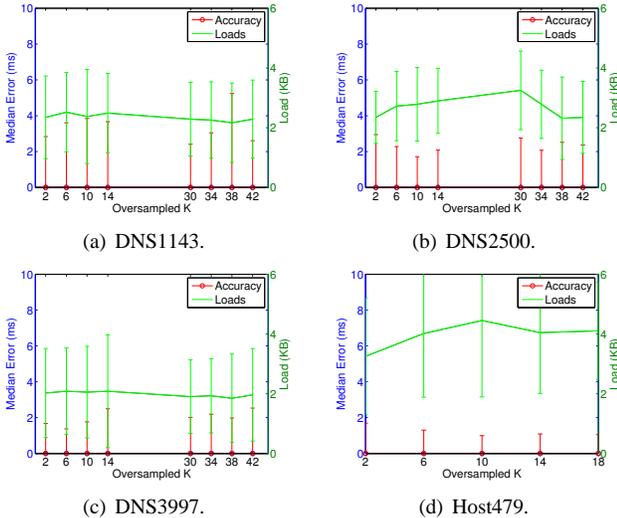


Fig. 18. Over-sampled number of neighbors.

7) *Returned Nodes For Next-Hop Probe m* : Fig. 19 plots the median errors and loads of HybridNN with increasing returned nodes for next-hop probes for HybridNN. For all data sets, HybridNN is accurate when the size of estimated nearest candidate neighbors for direct probes exceeds 2. Moreover, the loads of HybridNN increase slowly as the increment of relaxed probes. This is because we also add neighbors with higher uncertain coordinates, weakening the increased overhead of relaxed probes. Besides, the search process typically terminates at 3 to 5 hops as we found during experiments, therefore the measurement overhead is mostly bounded below 3 KB.

X. PLANETLAB EXPERIMENTS

We have implemented a prototype DNNs query system in Java using the asynchronous communication library. We implemented both HybridNN and Meridian. The core DNNs logic consists of around 5,000 lines of codes comprising three main modules: (1) prober module, which uses the kernel-level ping for delay measurements, to alleviate application level perturbations caused by high loads of PlanetLab nodes; (2) neighborhood management module, which finds and maintains neighbors on the concentric rings; (3) DNNs module, which utilizes the HybridNN or Meridian algorithm.

Our objective is to compare the accuracy and efficiency of DNNs queries with related nearest server location methods using real-world deployments. To that end, we choose 173 servers distributed globally on the PlanetLab as the service nodes. Then we select another 412 servers on the PlanetLab as the target machines. Our experiments last one week from 05-05-2011 to 12-05-2011.

We compare HybridNN with Meridian and iPlane [33]. We choose the same parameter configurations for HybridNN and Meridian as in the Simulation section (Sec IX-A). For iPlane, we query iPlane to obtain the delays between service nodes and target machines, then we compute the nearest service node for each target machine.

Besides, in order to compare the found nearest servers to the ground-truth nearest servers, we compute the ground-truth nearest servers using direct probes (denoted as *Direct*). Specifically, since pairwise delays between PlanetLab machines keep varying due to routing dynamics, we first use the median delay of any node pairs to summarize the long-term delay trend. Then we select the service node that has the lowest median delay value to the target.

A. Accuracy

First we compare the accuracy of different methods with the absolute error metric and the relative error metric defined in Sec IX-A. The results are shown in Fig 20(a) and (b).

HybridNN has significantly lower absolute errors and relative errors than Meridian. iPlane is similar with HybridNN, but incurs higher errors. The inaccuracy of iPlane is caused by the mismatch of the estimated routing paths and the real-world ones. The inaccuracy of Meridian shows that Meridian is easily trapped at local minimum far away from the optimal solutions.

On the other hand, HybridNN and iPlane are much accurate, which implies that hybridNN can avoid bad local minima in most cases. Nevertheless, HybridNN and iPlane also have around 3% of DNNS queries with relative errors above 10. We find that HybridNN incurs such high errors occur at the early stage, where nodes do not have enough neighbors in their concentric rings.

B. Completion Time

Next, we evaluate the completion time of individual DNNS queries for HybridNN and Meridian. Empirically, we have found that both HybridNN and Meridian complete DNNS queries within three search hops, which is consistent with the simulation results in Fig 12. However, the overall query time for DNNS searches depends on not only the number of search hops, but also the completion time of message exchanges and delay probes.

Fig 20(d) plots the distributions of query time of HybridNN and Meridian. Around 85% of the DNNS queries in HybridNN are similar with those of Meridian. Therefore, query time for HybridNN and Meridian are similar in most cases. However, around 20% of the queries take much large time to answer in Meridian, and 10% have query time larger than 15 seconds, while the hybrid measurement approach of HybridNN can avoid large query latencies.

C. Query Overhead

Next, to quantify the bandwidth overhead of the DNNS queries of HybridNN and Meridian, we define the load of a DNNS query as the total size of the transmitted packets during the DNNS process. We plot the CDFs of the loads for HybridNN and Meridian in Fig 20(d). The load of HybridNN is significantly lower than that of Meridian. In more than 95% of the cases the load of HybridNN is less than 2KBytes, while in more than 50% of the cases the load of Meridian is more than 10 KBytes, which is due to the large size of the candidate neighbor set for DNNS queries. Therefore, the delay estimation of HybridNN substantially reduces the measurement overhead.

D. Control Overhead

To measure the efficiency of HybridNN and Meridian. We collected the bandwidth overhead of the neighborhood management in HybridNN and Meridian for each service node every two minutes, as shown in Fig 20(e). The maintenance overhead of Meridian includes both the gossip process and the ring maintenance costs, while the maintenance of HybridNN includes the gossip messages, K nearest neighbor search messages and the K farthest neighbor search messages. The average maintenance overhead of HybridNN is 2 KBytes per

minute, and for Meridian is over 20 KBytes per minute. Since the time interval of ring maintenance for both HybridNN and Meridian is identical, the all-pair probes between nodes in the same ring is the main cause of the control overhead in Meridian. On the other hand, as HybridNN uses the coordinate distances to update the rings, it does not need to do all-pair probes between nodes in a ring.

XI. CONCLUSION AND FUTURE WORK

We have addressed the problem of designing an accurate and efficient DNNS algorithm in a comprehensive way. We first formulate the DNNS problem to account for both symmetric and asymmetric delay metrics for latency optimizations. Given the generalized delay metrics, we proposed to use the relaxed inframetric for modelling the delay space as a foundation for designing new DNNS algorithms with strong theoretical guarantees concerning search overhead and accuracy of the search results.

Next we apply all the insights gained to design a new DNNS algorithm called HybridNN. HybridNN locates nearest neighbors for any target using low bandwidth costs. For locating closer server to any target, HybridNN maximizes the diversity in the neighbor set, by discovering neighbors within each delay range through a light-weight neighbor sampling process. Next, in order to reduce the measurement costs of locating closer servers, HybridNN combines network coordinate based delay estimation and direct probes for fast and efficient nearest neighbor determination. Although the symmetric coordinate distances may deviate from the asymmetric delays, HybridNN is able to locate the nearest neighbor to the target at each search step, since we use direct probes to replace erroneous delay estimations. Finally, HybridNN terminates the search process conservatively in order to obtain better approximations of nearest neighbors. We confirmed the efficiency and effectiveness of HybridNN with extensive simulation and a prototype deployment on the PlanetLab. HybridNN can locate approximately closest neighbors quickly with low measurement costs.

As future work, we plan to continue two lines of research. First, currently we use the revised Vivaldi to estimate delays, which mismatches the asymmetric delay metric due to the symmetry of the coordinate distances. We plan to extend Vivaldi to asymmetric delay metrics. Second, we plan to study in-advance DNNS probing in order to hide the waiting time of on-demand DNNS queries for more practical latency-optimizations.

REFERENCES

- [1] R. Rodrigues and P. Druschel, "Peer-to-Peer Systems," *Commun. ACM*, vol. 53, no. 10, pp. 72–82, 2010.
- [2] Microsoft, "Office Live Workspace," <http://workspace.officelive.com/zh-hk/>, January 2011.
- [3] Google, "Google Maps," <http://maps.google.com/>, January 2011.
- [4] F. Agboma and A. Liotta, "QoE-aware QoS Management," in *Proc. of MoMM '08*, 2008, pp. 111–116.
- [5] A. G. Greenberg, J. R. Hamilton, D. A. Maltz, and P. Patel, "The Cost of a Cloud: Research Problems in Data Center Networks," *Computer Communication Review*, vol. 39, no. 1, pp. 68–73, 2009.
- [6] R. Krishnan, H. V. Madhyastha, S. Srinivasan, S. Jain, A. Krishnamurthy, T. E. Anderson, and J. Gao, "Moving Beyond End-to-End Path Information to Optimize CDN performance," in *Proc. of IMC'09*, pp. 190–201.

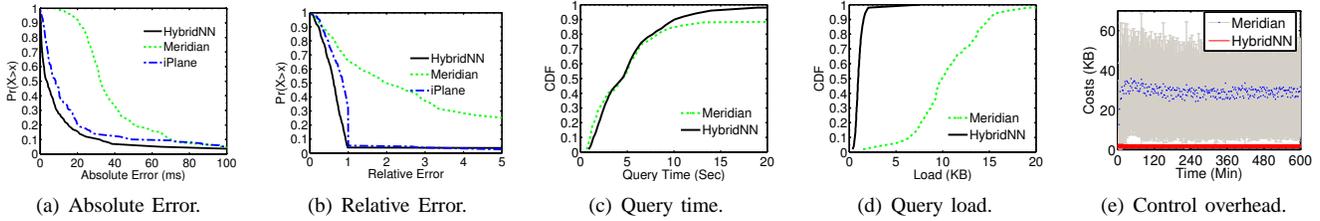


Fig. 20. Performance comparison on the PlanetLab.

[7] A.-J. Su, D. R. Choffnes, A. Kuzmanovic, and F. E. Bustamante, "Drafting behind Akamai (travelocity-based detouring)," in *Proc. of SIGCOMM'06*.

[8] M. J. Freedman, E. Freudenthal, and D. Mazières, "Democratizing content publication with coral," in *Proc. of NSDI'04*, pp. 18–18.

[9] M. J. Freedman, K. Lakshminarayanan, and D. Mazières, "OASIS: Anycast for Any Service," in *Proc. of NSDI'06*.

[10] P. Wendell, J. W. Jiang, M. J. Freedman, and J. Rexford, "DONAR: Decentralized Server Selection for Cloud Services," in *Proc. of SIGCOMM'10*, pp. 231–242.

[11] J. D. Guyton, J. D. Guyton, and M. F. Schwartz, "Locating Nearby Copies of Replicated Internet Servers," in *Proc. of SIGCOMM '95*, pp. 288–298.

[12] M. Costa, M. Castro, A. I. T. Rowstron, and P. B. Key, "PIC: Practical Internet Coordinates for Distance Estimation," in *Proc. of ICDCS'04*, pp. 178–187.

[13] B. Wong, A. Slivkins, and E. G. Sirer, "Meridian: a Lightweight Network Location Service Without Virtual Coordinates," in *Proc. of SIGCOMM'05*, pp. 85–96.

[14] V. Vishnumurthy and P. Francis, "On the Difficulty of Finding the Nearest Peer in P2P Systems," in *Proc. of IMC'08*, pp. 9–14.

[15] C. Lumezanu, R. Baden, N. Spring, and B. Bhattacharjee, "Triangle Inequality Variations in the Internet," in *Proc. of IMC '09*, pp. 177–183.

[16] A. Pathak, H. Pucha, Y. Zhang, Y. C. Hu, and Z. M. Mao, "A measurement study of internet delay asymmetry," in *Proc. of PAM'08*, pp. 182–191.

[17] S. Shalunov, B. Teitelbaum, A. Karp, J. Boote, and M. Zekauskas, "A One-way Active Measurement Protocol (OWAMP)," RFC 4656 (Proposed Standard), 2006. [Online]. Available: <http://www.ietf.org/rfc/rfc4656.txt>

[18] E. Katz-Bassett, H. V. Madhyastha, V. K. Adhikari, C. Scott, J. Sherry, P. van Wesep, T. E. Anderson, and A. Krishnamurthy, "Reverse traceroute," in *Proc. of NSDI'10*, pp. 219–234.

[19] V. Paxson, "End-to-end internet packet dynamics," in *Proc. of SIGCOMM '97*, 1997, pp. 139–152.

[20] G. Wang, B. Zhang, and T. S. E. Ng, "Towards Network Triangle Inequality Violation Aware Distributed Systems," in *Proc. of IMC'07*, pp. 175–188.

[21] F. Dabek, R. Cox, M. F. Kaashoek, and R. Morris, "Vivaldi: a Decentralized Network Coordinate System," in *Proc. of SIGCOMM'04*, pp. 15–26.

[22] S. Ratnasamy, M. Handley, R. Karp, and S. Shenker, "Topologically-aware Overlay Construction and Server Selection," in *Proc. of INFOCOM'02*, pp. 1190 – 1199 vol.3.

[23] M. Waldvogel and R. Rinaldi, "Efficient Topology-Aware Overlay Network," in *Proc. of Hotnets-I*, 2002.

[24] G. R. Hjaltason and H. Samet, "Index-driven Similarity Search in Metric Spaces (Survey Article)," *ACM Trans. Database Syst.*, vol. 28, pp. 517–580, 2003.

[25] K. L. Clarkson, "Nearest-Neighbor Searching and Metric Space Dimensions," in *Nearest-Neighbor Methods for Learning and Vision: Theory and Practice*, G. Shakhnarovich, T. Darrell, and P. Indyk, Eds. MIT Press, 2006, pp. 15–59.

[26] E. Chávez, G. Navarro, R. Baeza-Yates, and J. L. Marroquín, "Searching in Metric Spaces," *ACM Comput. Surv.*, vol. 33, pp. 273–321, 2001.

[27] P. Indyk. (2004) Nearest Neighbors In High-Dimensional Spaces. [Online]. Available: <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.10.3826>

[28] S. M. Hotz, "Routing information organization to support scalable interdomain routing with heterogeneous path requirements," PhD Thesis, Computer Science Department, University of Southern California, Los Angeles, California, 1994.

[29] R. L. Carter and M. E. Crovella, "Server selection using dynamic path characterization in wide-area networks," in *Proc. of INFOCOM '97*, pp. 1014–.

[30] —, "On the network impact of dynamic server selection," *Computer Networks*, vol. 31, no. 23-24, pp. 2529 – 2558, 1999.

[31] P. Sharma, Z. Xu, S. Banerjee, and S.-J. Lee, "Estimating network proximity and latency," *Computer Communication Review*, vol. 36, no. 3, pp. 39–50, 2006.

[32] A.-J. Su, D. Choffnes, F. E. Bustamante, and A. Kuzmanovic, "Relative network positioning via cdn redirections," in *Proc. of ICDCS '08*, pp. 377–386.

[33] H. V. Madhyastha, E. Katz-Bassett, T. E. Anderson, A. Krishnamurthy, and A. Venkataramani, "iPlane Nano: Path Prediction for Peer-to-Peer Applications," in *Proc. of NSDI'09*, pp. 137–152.

[34] H. V. Madhyastha, T. Isdal, M. Piatek, C. Dixon, T. E. Anderson, A. Krishnamurthy, and A. Venkataramani, "iPlane: An Information Plane for Distributed Services," in *Proc. of OSDI'06*, pp. 367–380.

[35] S. Banerjee, C. Kommareddy, and B. Bhattacharjee, "Scalable Peer Finding on the Internet," in *Proc. of Global Internet Symposium 2002*.

[36] Y. Fu, Y. Wang, and E. Biersack, "HybridNN: Supporting Network Location Service on Generalized Delay Metrics for Latency Sensitive Applications," Eurecom, Technical Report 1, January 2011.

[37] B. Zhang, T. S. E. Ng, A. Nandi, R. H. Riedi, P. Druschel, and G. Wang, "Measurement-based analysis, modeling, and synthesis of the internet delay space," *IEEE/ACM Trans. Netw.*, vol. 18, no. 1, pp. 229–242, 2010.

[38] K. P. Gummadi, S. Saroiu, and S. D. Gribble, "King: Estimating Latency Between Arbitrary Internet End Hosts," in *Proc. of IMW '02*, pp. 5–18.

[39] D. R. Choffnes, M. Sanchez, and F. E. Bustamante, "Network Positioning from the Edge - An Empirical Study of the Effectiveness of Network Positioning in P2P Systems," in *Proc. of INFOCOM'10*, pp. 291–295.

[40] P2PSim, "The P2PSim Project," <http://pdos.csail.mit.edu/p2psim/kingdata/>, October 2010.

[41] Y. He, M. Faloutsos, S. Krishnamurthy, and B. Huffaker, "On routing asymmetry in the internet," in *Proc. of GLOBECOM '05*.

[42] D. R. Choffnes and F. E. Bustamante, "Pitfalls for testbed evaluations of internet systems," *SIGCOMM Comput. Commun. Rev.*, vol. 40, pp. 43–50, April 2010.

[43] P. Fraigniaud, E. Lebar, and L. Viennot, "The Inframetric Model for the Internet," in *Proc. of INFOCOM'08*, pp. 1085–1093.

[44] D. R. Karger and M. Ruhl, "Finding Nearest Neighbors in Growth-restricted Metrics," in *Proc. of STOC '02*, 2002, pp. 741–750.

[45] J. Ledlie, P. Gardner, and M. I. Seltzer, "Network Coordinates in the Wild," in *Proc. of NSDI'07*.

[46] Y. Zhang and N. G. Duffield, "On the Constancy of Internet Path Properties," in *Proc. of IMW'01*, pp. 197–211.

APPENDIX

Lemma VI.1: Given a ρ -inframetric with growth $\gamma_g \geq 1$, for any $x \geq \rho$, $r > 0$ and any node P , the volume of a ball $B_P(r)$ is at most x^α smaller than that of the ball $B_P(xr)$, where $\log_\rho \gamma_g \leq \alpha \leq 2 \log_\rho \gamma_g$.

Proof: First, according to the definition of the growth, it follows:

$$|B_P(xr)| \leq \gamma_g \left| B_P\left(\frac{x}{\rho}r\right) \right|$$

Then, by recursively calling $\lceil \log_\rho x \rceil$ times the growth definition, until $\frac{x}{\rho^{\lceil \log_\rho x \rceil}} < 1$, then

$$\begin{aligned} |B_P(xr)| &\leq \gamma_g^{\lceil \log_\rho x \rceil} |B_P(r)| = x^{\log_x \gamma_g^{\lceil \log_\rho x \rceil}} |B_P(r)| \\ &= x^\alpha |B_P(r)|, \alpha = \log_x \gamma_g \times \lceil \log_\rho x \rceil \end{aligned}$$

Therefore, by the definition of the ceiling function, we can calculate the lower bound of α as:

$$\alpha \geq \log_x \gamma_g \times \log_\rho x = \log_\rho \gamma_g$$

On the other hand, due to $x \geq \rho$, $\gamma_g > 1$, we get

$$\log_\rho \gamma_g = \frac{\log \gamma_g}{\log \rho} \geq \frac{\log \gamma_g}{\log x} = \log_x \gamma_g$$

thus we can compute the upper bound of α as:

$$\begin{aligned} \alpha &\leq \log_x \gamma_g \times (\log_\rho x + 1) \\ &= \log_\rho \gamma_g + \log_x \gamma_g \\ &\leq \log_\rho \gamma_g + \log_\rho \gamma_g \\ &= 2\log_\rho \gamma_g \end{aligned}$$

this concludes the proof. \blacksquare

Lemma VI.2: (Sandwich lemma) For any pair of node p and q , and $d_{pq} \leq r$, then

$$B_q(r) \subseteq B_p(\rho r) \subseteq B_q(\rho^2 r)$$

Proof: (1) For any node i satisfying $d_{qi} \leq r$, i.e., $i \in B_q(r)$, by the definition of the inframetric model, $d_{pi} \leq \rho \max\{d_{pq}, d_{qi}\} \leq \rho r$, thus $i \in B_p(\rho r)$, that is,

$$B_q(r) \subseteq B_p(\rho r)$$

(2) For any node j satisfying $j \in B_p(\rho r)$, by the definition of the inframetric model, it follows

$$d_{qj} \leq \rho \{d_{pq}, d_{pj}\} \leq \rho^2 r$$

Summing up (1) and (2) conclude the proof. \blacksquare

Theorem VI.3: (Sampling efficiency in the growth dimension) For a ρ -inframetric model with growth $\gamma_g \geq 1$, for a service node P , and a DNNS target T satisfying $d_{PT} \leq r$, when selecting $3\left(\frac{\rho^2}{\beta}\right)^\alpha$ nodes uniformly at random from $B_P(\rho r)$ with replacement, with probability of at least 95%, one of these nodes will lie in $B_T(\beta r)$, where $\log_\rho \gamma_g \leq \alpha \leq 2\log_\rho \gamma_g$ and $\beta < 1$.

Proof: since $B_T(\beta r) \subset B_T(r) \subseteq B_P(\rho r)$ by the sandwich lemma VI.2, all nodes covered by $B_T(\beta r)$ are also covered by $B_P(\rho r)$. Therefore, we only need to sample enough nodes in $B_P(\rho r)$ in order to sample a node located in $B_T(\beta r)$.

Furthermore, for the pair of nodes P and T satisfying $d_{PT} \leq r$, it follows

$$|B_P(\rho r)| \leq |B_T(\rho^2 r)| = \left| B_T\left(\frac{\rho^2}{\beta} \beta r\right) \right|$$

Since we know $\rho > 1$, then $\frac{\rho^2}{\beta} > \rho^2 > \rho$, therefore the preconditions of lemma VI.1 hold, by lemma VI.1, we can

show the relation between the ball $B_P(\rho r)$ and the ball $B_T(\beta r)$ where $\beta < 1$,

$$|B_P(\rho r)| \leq \left| B_T\left(\frac{\rho^2}{\beta} \beta r\right) \right| \leq \left(\frac{\rho^2}{\beta}\right)^\alpha |B_T(\beta r)|$$

where $\log_\rho \gamma_g \leq \alpha \leq 2\log_\rho \gamma_g$. Therefore, the probability of uniformly sampling a node from $B_P(\rho r)$ which lies in the ball $B_T(\beta r)$ is:

$$\frac{|B_T(\beta r)|}{|B_P(\rho r)|} \geq \frac{|B_T(\beta r)|}{\left(\frac{\rho^2}{\beta}\right)^\alpha |B_T(\beta r)|} = \frac{1}{\left(\frac{\rho^2}{\beta}\right)^\alpha}$$

Consequently, the probability that $3\left(\frac{\rho^2}{\beta}\right)^\alpha$ samples are not in the ball $B_T(\beta r)$ is at most

$$\left(1 - \frac{1}{\left(\frac{\rho^2}{\beta}\right)^\alpha}\right)^{3\left(\frac{\rho^2}{\beta}\right)^\alpha} \leq \left(\frac{1}{e}\right)^3 \leq 0.05$$

Thus, with probability more than 95% we succeed in locating a node lying in the ball $B_T(\beta r)$ with $3\left(\frac{\rho^2}{\beta}\right)^\alpha$ samples. \blacksquare

Corollary A.1. *For a relaxed inframetric model with growth γ_g , according to the DNNS process in Definition VI.4, the found nearest neighbor is a $\frac{1}{\beta}$ -approximation, and the number of search steps is smaller than $\log_{\frac{1}{\beta}} \Delta$, where Δ is the ratio of the maximum delay to the minimum delay of all pairwise delays.*

Proof: If a DNNS request is forwarded from node P to node Q , the progress is said to be $\frac{d_{PT}}{d_{QT}}$. According to the DNNS search process, by Theorem VI.3, the progress is at least $\frac{1}{\beta}$ at every node P , therefore in at most $\log_{\frac{1}{\beta}} \Delta$ steps, we reach some node v satisfying $d_{vT} < \frac{1}{\beta} d_*$, which terminates the DNNS query process as we can not find suitable next-hop neighbors, where d_* is the minimum delay to target T . Therefore, the found nearest neighbor v is $\frac{1}{\beta}$ -approximation. \blacksquare