

Vivaldi

Practical, Distributed Internet Coordinates

Russ Cox, Frank Dabek,
Frans Kaashoek, Robert Morris,
and many others

rsc@mit.edu

Computer Science and Artificial Intelligence Lab
Massachusetts Institute of Technology

~~January 11, 2005~~

February 7, 2005

at IDA-CCS, Bowie, Maryland

Circa 2001: file sharing networks extremely popular

- harness huge numbers of machines
- but don't scale well (at least asymptotically)

Academics study how to build scalable peer-to-peer systems.

- Chord, Pastry, OceanStore, ...

How can we make them perform better in practice?

- predict round trip times with Vivaldi

P2P circa 2001: an exciting social development

- Internet users cooperating to share, for example, music files.
 - Napster, Gnutella, Morpheus, KaZaA, ...
- Lots of attention from the popular press (and the RIAA!)
 - “The ultimate form of democracy on the Internet.”
 - “The ultimate threat to copyright protection on the Internet.”
- Not much new technologically.

What is a P2P System?

- System without any central servers.
 - Every node is a server
 - No particular node is vital to the network
 - Nodes all have the same functionality
- Huge number of nodes, many node failures
- Enabled by technology improvements

P2P: useful for building reliable services?

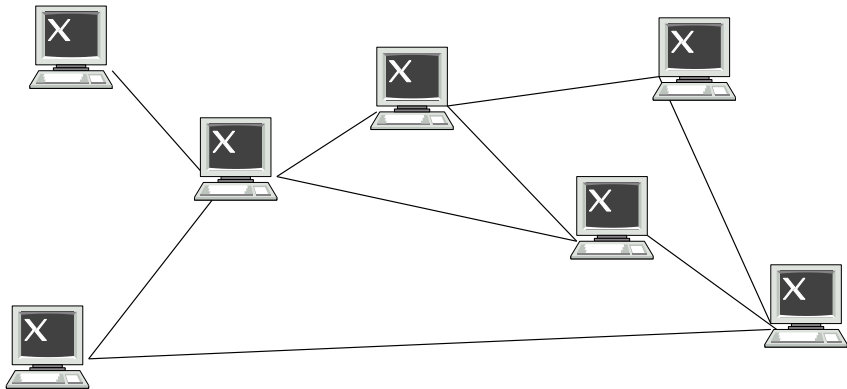
- Many critical services use the Internet.
 - Hospitals, some government agencies, etc.
 - (but not other Agencies)
- Non-Internet systems exist at large scales too.
 - corporations, government agencies, etc.
- Can we build large-scale robust distributed services?
 - Node and communication failures
 - Load fluctuations (e.g., flash crowds)
 - Attacks (including DDoS)

The promise of P2P computing

- Reliability: no central point of failure.
 - Many replicas
 - Geographic distribution
- High capacity through parallelism
 - Many disks
 - Many network connections
 - Many CPUs
- Automatic configuration
- Useful in public and proprietary settings

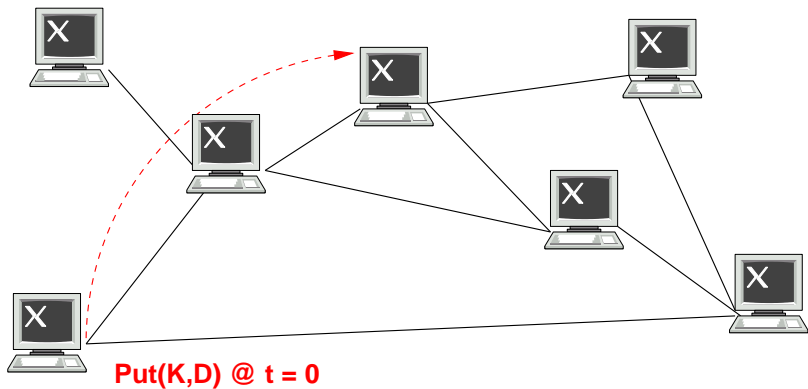
Distributed Hash Table

- Building block for “principled” peer-to-peer systems.

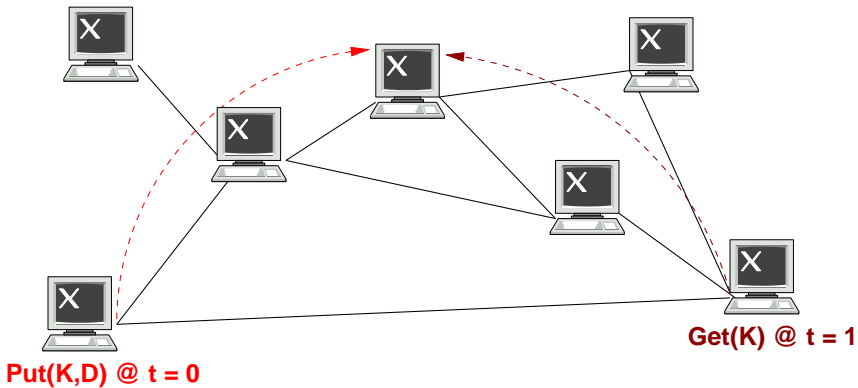


- DHTs provide location service: map key to a set of machines.

DHT as Location and Rendezvous Service



DHT as Location and Rendezvous Service



DHT as Location and Rendezvous Service

Many useful primitives can be and have been built on top of this rendezvous.

- File sharing
- Web cache
- Archival/Backup storage
- Censor-resistant stores
- DB query and indexing
- Event notification
- Naming systems
- Communication primitives

We will concentrate on storage.

- but first, how to implement DHT?

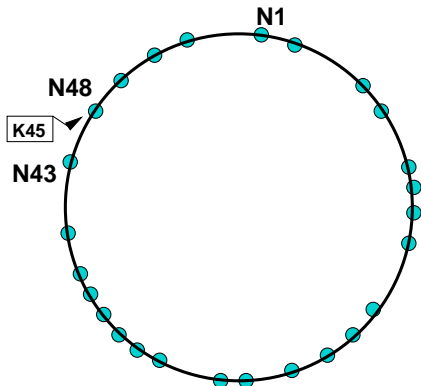
Chord: One Peer-to-Peer System

Developed at MIT in 2000.

Recipe differs from others, but flavor is same.

- Assign nodes and keys big identifiers (SHA1 hashes).
- Traverse identifier space during lookup.
- Provide lookup first; layer storage on top.

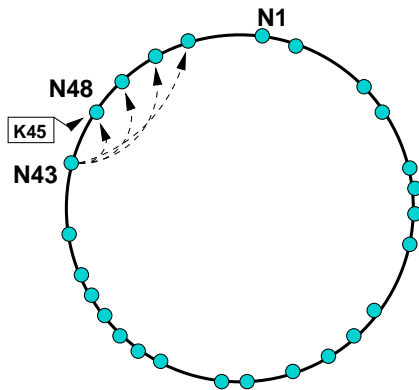
Chord: Successor Defines Ownership



Identifier space arranged in a *logical* ring.

- A key's *successor* is first node clockwise after key on ring.
- That node is the owner of the key.

Chord: Successors Ensure Correctness



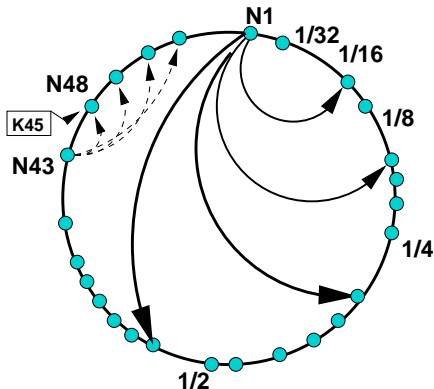
Nodes maintain list of r successors.

- On failure of node n , successor takes over keys.
- Nodes ping predecessors regularly to detect failure.

Can traverse ring by following successors.

- a bit slow if there are millions of nodes.

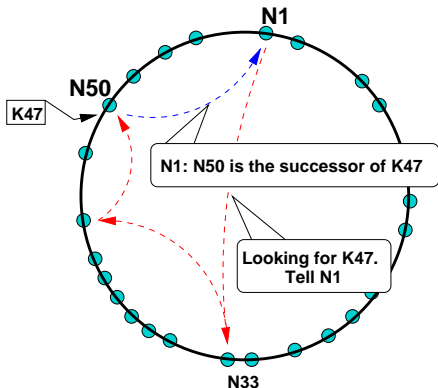
Chord: Fingers Speed Lookups



Nodes maintain pointers to larger hops around the ring.

- Like using skip lists in a linked list.
- $O(\log N)$ exponentially-distributed fingers.

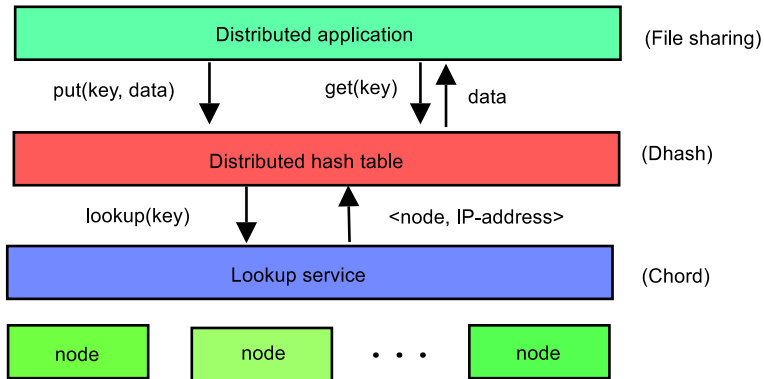
Chord: Lookup Requires $\log N$ Hops



Each hop halves the distance to the key.

DHash: a Storage Infrastructure

Now we have lookup; add storage (and then apps).



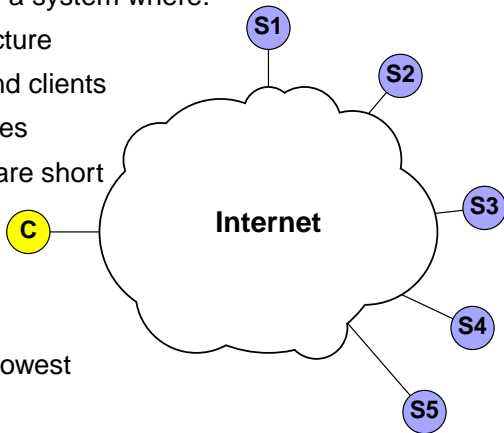
To fetch a block in DHash:

- Block replicated on r successors of key
- Use Chord (DHT) to find servers responsible for block
- Download block from one of those servers. Which?

Problem: Predicting Round Trip Times on Internet

Example: server selection in a system where:

- no centralized infrastructure
- nodes act as servers and clients
- many thousands of nodes
- exchanges with server are short
- server choice changes for each exchange



Want to choose server with lowest round trip time to client.

How?

- Can avoid predictions, wasting time or bandwidth:
 - measure RTT on demand
 - measure RTT in advance
 - talk to multiple servers at once
- Can predict using synthetic coordinates as in GNP (Infocom 2002).

Synthetic Coordinates with GNP

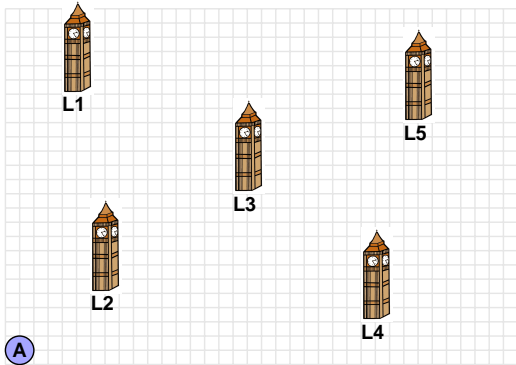
GNP assigns Euclidean coordinates to nodes such that coordinate distance predicts round trip time.

Synthetic Coordinates with GNP

GNP assigns Euclidean coordinates to nodes such that coordinate distance predicts round trip time.

Node A pings landmarks to compute its own position.

	Coord	Dist
L1	(40,320)	
L2	(60,180)	
L3	(160,250)	
L4	(250,160)	
L5	(280,300)	

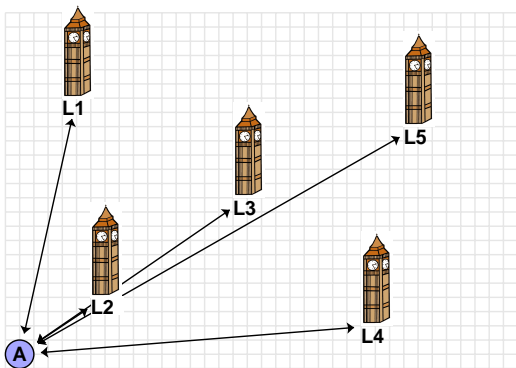


Synthetic Coordinates with GNP

GNP assigns Euclidean coordinates to nodes such that coordinate distance predicts round trip time.

Node A pings landmarks to compute its own position.

	Coord	Dist
L1	(40,320)	117 ms
L2	(60,180)	201 ms
L3	(160,250)	110 ms
L4	(250,160)	223 ms
L5	(280,300)	143 ms

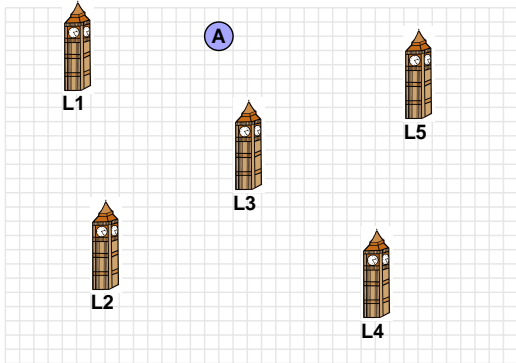


Synthetic Coordinates with GNP

GNP assigns Euclidean coordinates to nodes such that coordinate distance predicts round trip time.

Node A pings landmarks to compute its own position.

	Coord	Dist
L1	(40,320)	117 ms
L2	(60,180)	201 ms
L3	(160,250)	110 ms
L4	(250,160)	223 ms
L5	(280,300)	143 ms

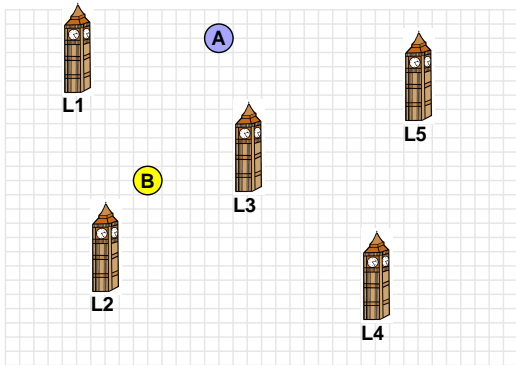


Synthetic Coordinates with GNP

GNP assigns Euclidean coordinates to nodes such that coordinate distance predicts round trip time.

Node A pings landmarks to compute its own position.

Node B does the same.



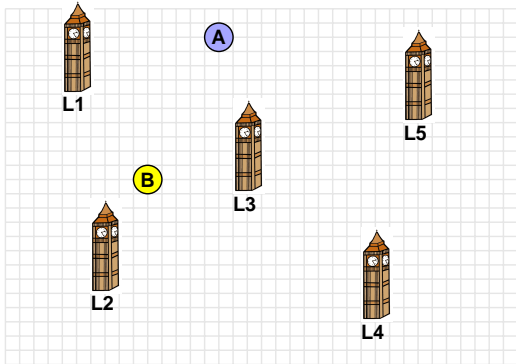
Synthetic Coordinates with GNP

GNP assigns Euclidean coordinates to nodes such that coordinate distance predicts round trip time.

Node A pings landmarks to compute its own position.

Node B does the same.

RTT between A and B is predicted by the distance between their coordinates, *without direct measurement*.

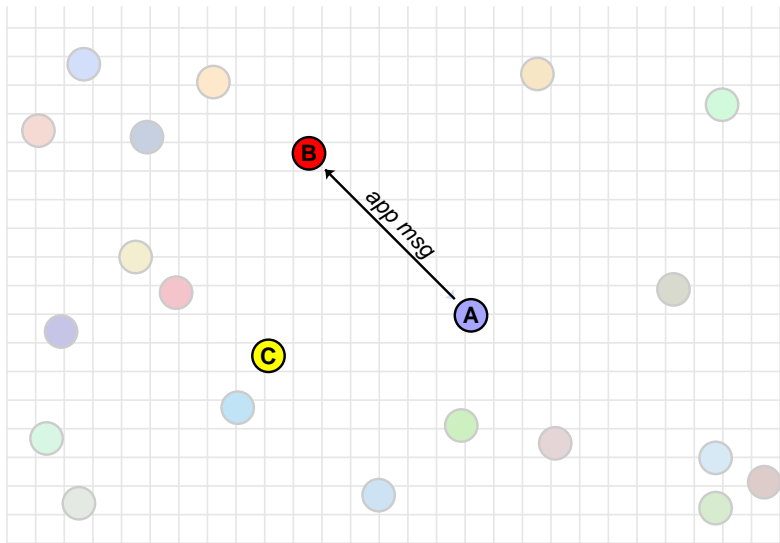


Vivaldi is a decentralized method for computing synthetic coordinates

- Piggyback on application traffic
- Node updates its own coordinates in response to sample
- Each node need only contact a small fraction of the other nodes

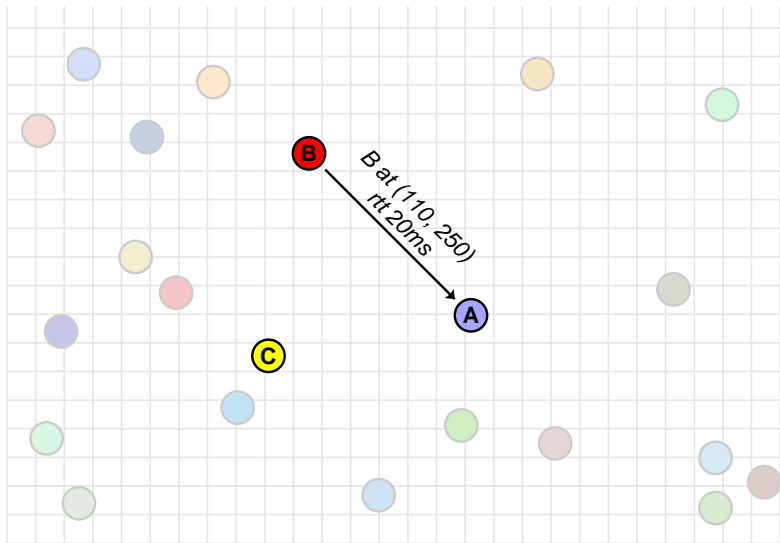
Vivaldi Example

Follow node A through a sequence of communications.



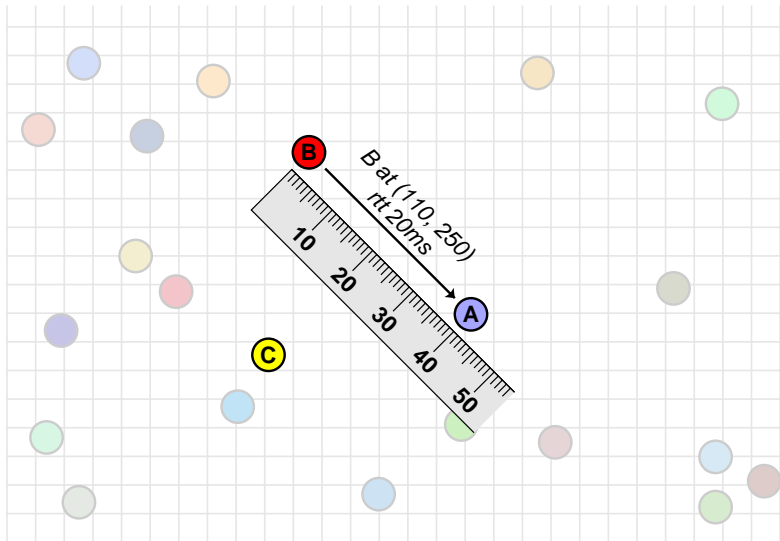
Vivaldi Example

A obtains B's coordinates, RTT.



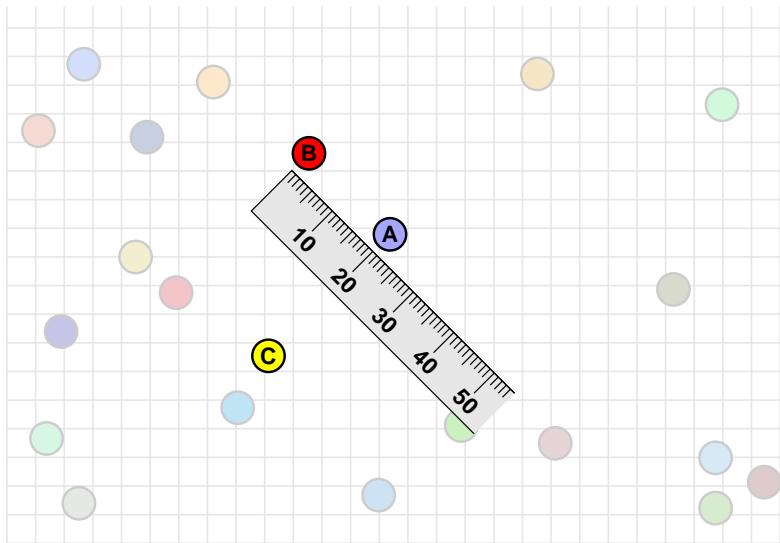
Vivaldi Example

A computes distance to B in coordinate space.



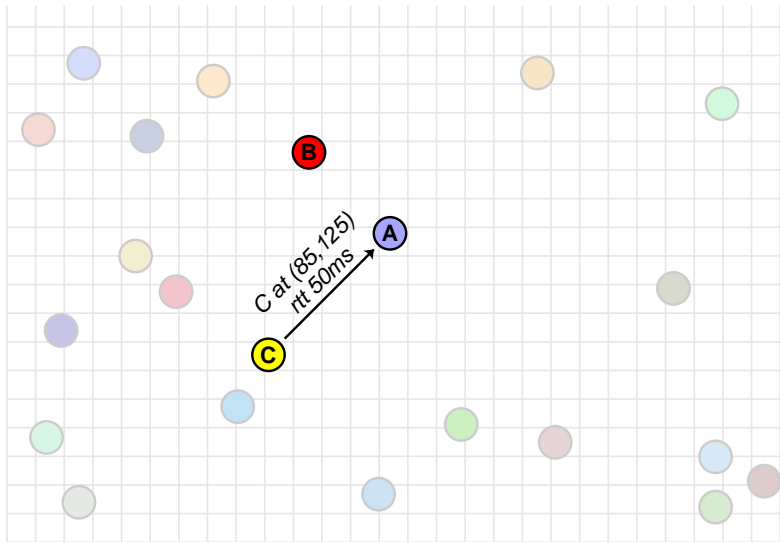
Vivaldi Example

A adjusts coordinates so distance matches actual RTT.



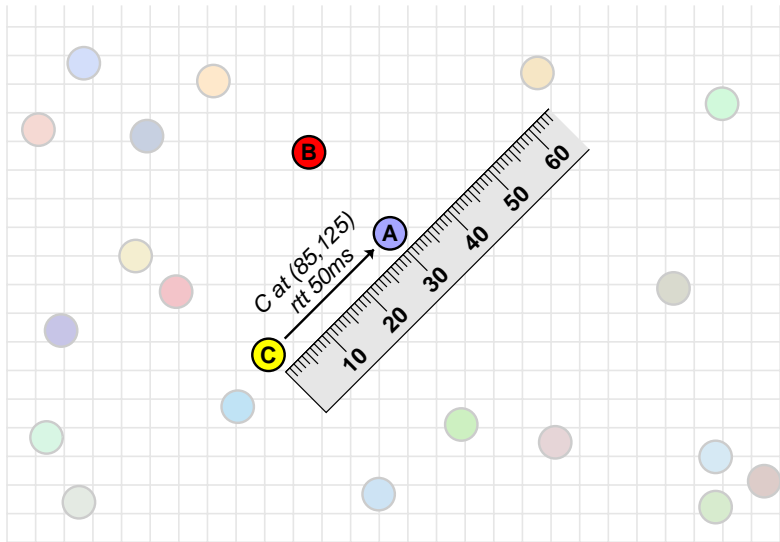
Vivaldi Example

A obtains C's coordinates, RTT.



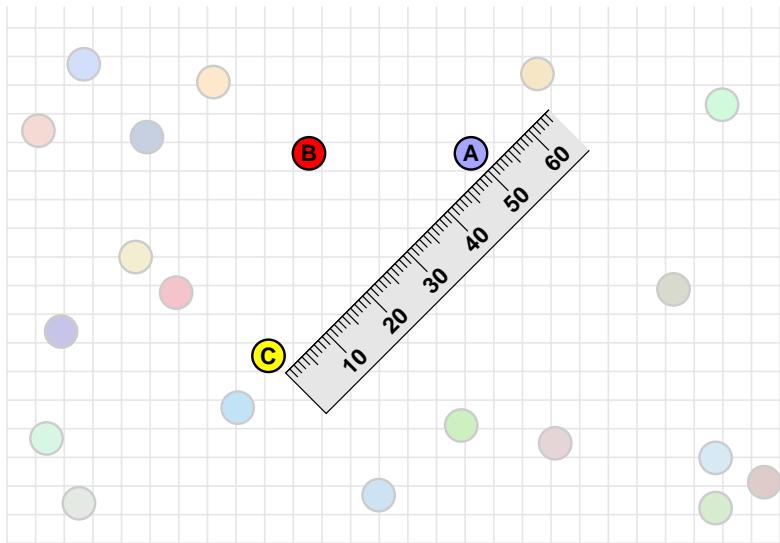
Vivaldi Example

A computes distance to C in coordinate space.



Vivaldi Example

A adjusts coordinates so distance matches actual RTT.
(Now A is wrong distance from B.)



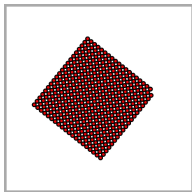
Challenges of Decentralization

Without centralized control, must consider:

- will the system converge to an accurate coordinate set?
- how long will the system take to converge?
- will the system be disturbed by new nodes joining the system?

Tuning Vivaldi: Convergence

- Run Vivaldi on round trip times derived from grid.



- As described, algorithm never converges.



- To cause convergence, damp motion.

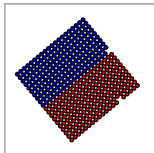


- To speed convergence, vary damping with estimate of prediction accuracy.



Tuning Vivaldi: Naive Newcomers

- Run Vivaldi on round trip times derived from grid.
Blue nodes start first, stabilize.
Red nodes join the system.



- High-accuracy nodes are displaced by new, low-accuracy nodes joining the system
- To avoid this, vary damping with ratio of local node's accuracy and sampled node's accuracy.



Vivaldi Algorithm

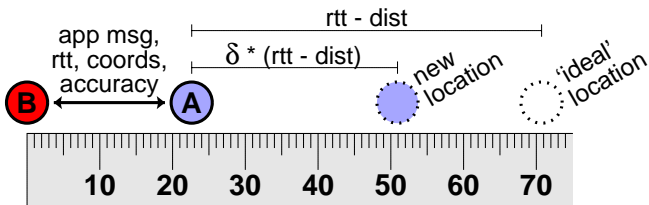
Given the coordinates, round trip time, and accuracy estimate of a node:

- Update local accuracy estimate.
- Compute 'ideal' location.
- Compute damping constant δ using local and remote accuracy estimates.
- Move δ of the way toward the "ideal" location.

Vivaldi Algorithm

Given the coordinates, round trip time, and accuracy estimate of a node:

- Update local accuracy estimate.
- Compute 'ideal' location.
- Compute damping constant δ using local and remote accuracy estimates.
- Move δ of the way toward the "ideal" location.



Evaluating Synthetic Coordinates on the Internet

- Cannot evaluate by comparing to “correct” coordinate set.
- Evaluate *predictions* made using a coordinate set.
- Predictions of Internet will never be perfect.
 - violations of triangle inequality, ...

Evaluating Vivaldi on the Internet

- How accurate are Vivaldi's predictions?
- How quickly does Vivaldi converge to a coordinate set?
- How quickly can Vivaldi adapt to network changes?
- How does choice of coordinate space affect error?
- How does Vivaldi work in real-world apps?

Use simulator seeded with real Internet measurements.

- pairwise RTTs for 192 PlanetLab nodes
- use RTT matrix as input to simulator
- run various algorithms on simulator

Each Vivaldi node queries others as fast as it can

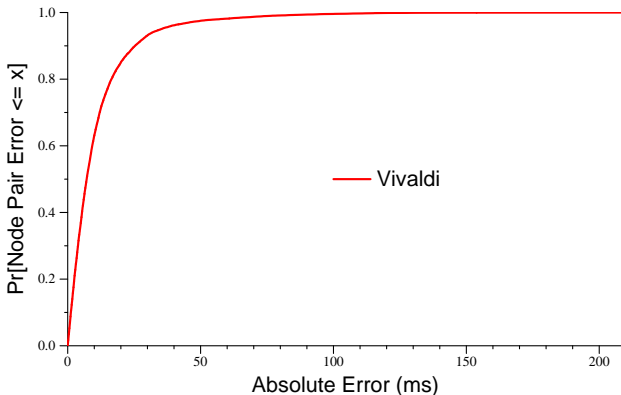
- one message outstanding at a time
- each node has a small fixed neighbor set

Vivaldi's Absolute Prediction Error

Look at distribution of absolute prediction error, defined as

$$|\text{actual RTT} - \text{predicted RTT}|,$$

over all node pairs in the system.

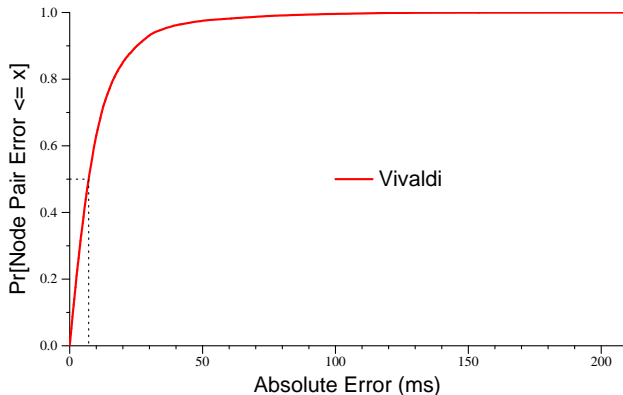


Vivaldi's Absolute Prediction Error

Look at distribution of absolute prediction error, defined as

$$|\text{actual RTT} - \text{predicted RTT}|,$$

over all node pairs in the system.

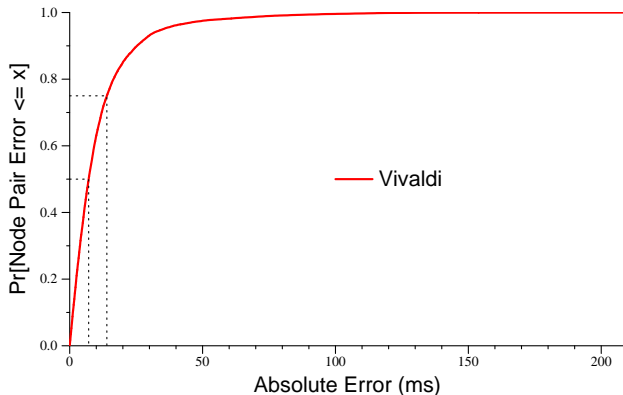


Vivaldi's Absolute Prediction Error

Look at distribution of absolute prediction error, defined as

$$|\text{actual RTT} - \text{predicted RTT}|,$$

over all node pairs in the system.

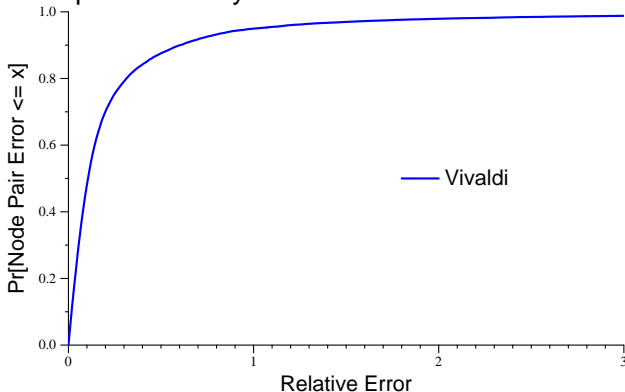


Vivaldi's Relative Prediction Error

Look at relative error, defined as

$$\frac{|\text{actual RTT} - \text{predicted RTT}|}{\min(\text{actual RTT}, \text{predicted RTT})},$$

over all node pairs in the system.

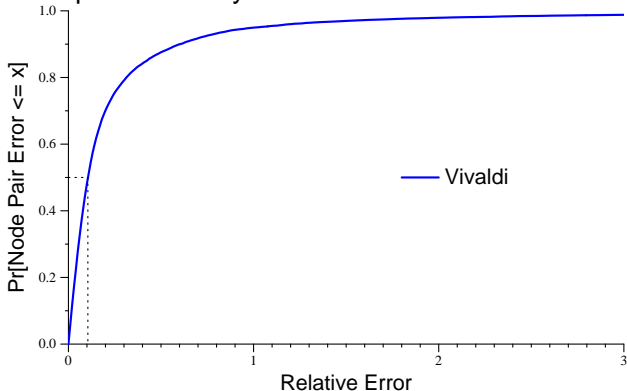


Vivaldi's Relative Prediction Error

Look at relative error, defined as

$$\frac{|\text{actual RTT} - \text{predicted RTT}|}{\min(\text{actual RTT}, \text{predicted RTT})},$$

over all node pairs in the system.

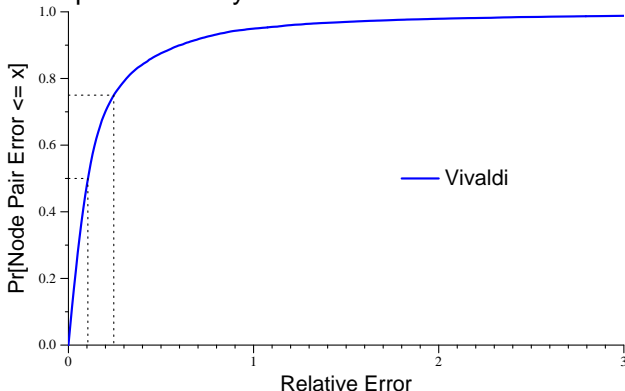


Vivaldi's Relative Prediction Error

Look at relative error, defined as

$$\frac{|\text{actual RTT} - \text{predicted RTT}|}{\min(\text{actual RTT}, \text{predicted RTT})},$$

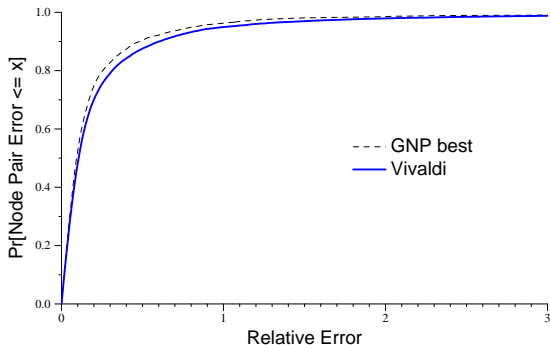
over all node pairs in the system.



Vivaldi Compared to GNP on Relative Error

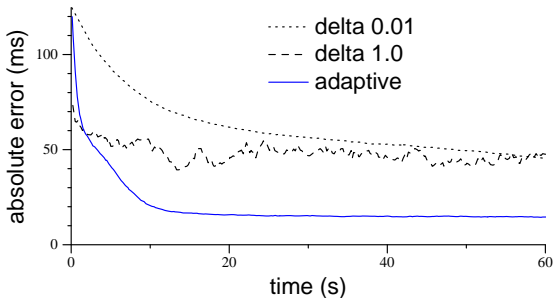
Compare to GNP's predictions.

- GNP sensitive to landmark choice. Use best of 64 random landmark sets.



Vivaldi's Convergence Time

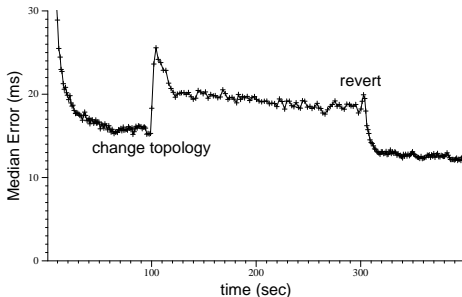
- Depends on choice of δ , the damping constant.



- Using adaptive δ , Vivaldi converges in under 20 seconds (60 measurements per node).

Vivaldi's Time to Adapt to Network Changes

- Vivaldi nodes are always adjusting their coordinates.
- Test adapting speed with synthetic topology change: lengthen one link by factor of ten.



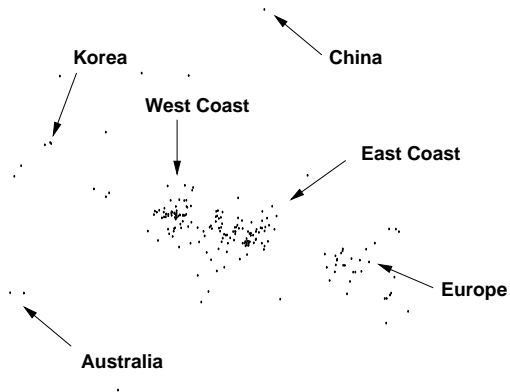
- Vivaldi adapts in about twenty seconds.

Other Coordinate Spaces

- *A priori*, it's not clear why *any* coordinate system should fit the Internet well.
- GNP showed that Euclidean coordinates work well.
- Why do they work?
- Are there better coordinate systems?
 - Obvious other candidates: globe, 3D, 4D, ...

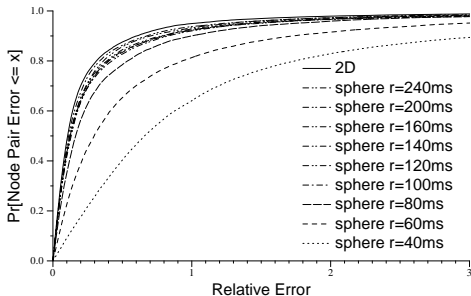
Vivaldi's 2D Assignment for PlanetLab

Placement in 2D mirrors physical geography.



Globe Coordinates vs. 2D Euclidean

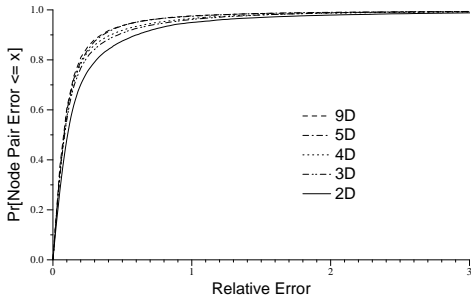
- Globe coordinates (latitude, longitude) place nodes on surface of a sphere.
- Great circle distance between two nodes on the sphere depends on radius.



- Coordinate sets are using one part of the sphere as a rough approximation to a 2D plane.

Higher Euclidean Dimensions

- If two are good, more should be better.



- Why are they better?

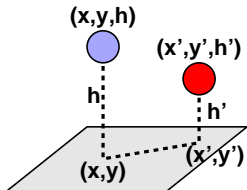
Higher Euclidean Dimensions Explained

- In 2D, some nodes need to be farther away from all others.
- In 3D, these “hard-to-place” nodes can move up or down from the 2D plane to get away from everyone.
- Each new dimension adds an independent direction.
- Accommodates per-node overhead: server load, access links.

Problem: how can we accommodate “hard-to-place” nodes without an arbitrary number of dimensions?

Height Vectors

- Give “hard-to-place” nodes their own way to get away from everyone.
- Height vectors place nodes at some height above a 2D transit plane.
- Directly models per-node overhead.

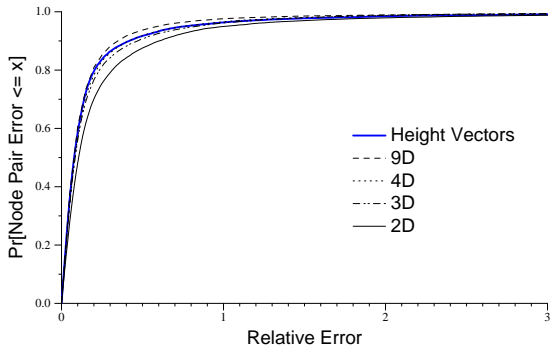


- Distance from (x, y, h) to (x', y', h') is

$$h + \sqrt{(x - x')^2 + (y - y')^2} + h'.$$

Height Vectors Work Well

- Height Vectors outperform 2- and 3-D Euclidean.



- Works to view Internet as geographically-accurate core with access links attached.

Evaluating Vivaldi on Chord

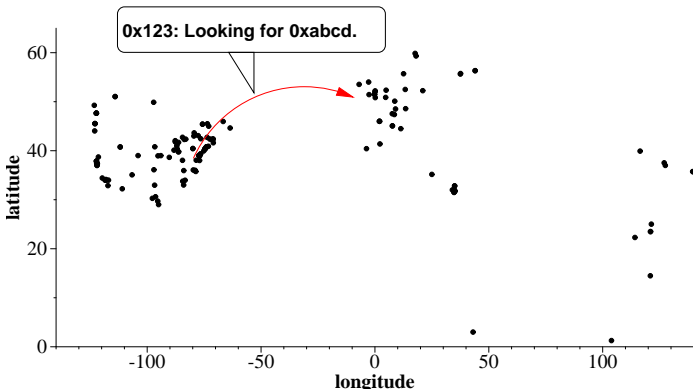
Chord performance is $O(\log n)$ but sluggish.

Can Vivaldi help?

Watching Traffic

Lookups move randomly in the network.

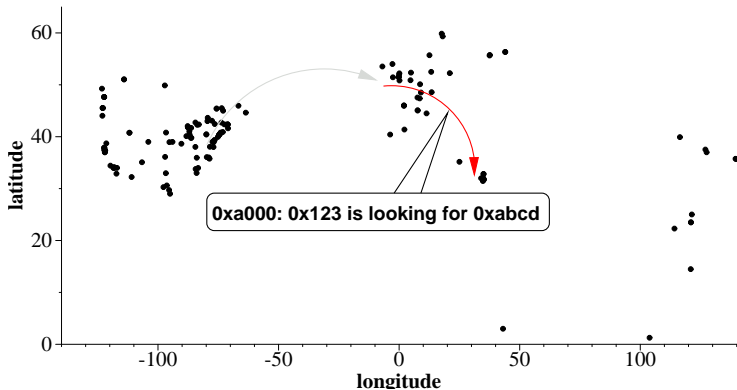
- Node IDs, block keys are assigned (pseudo-)randomly.



Watching Traffic

Lookups move randomly in the network.

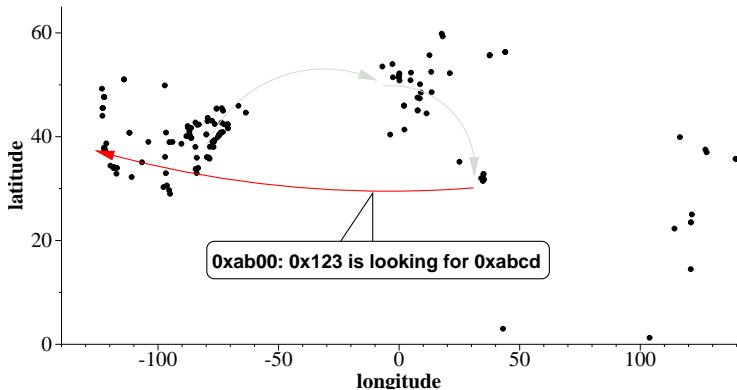
- Node IDs, block keys are assigned (pseudo-)randomly.



Watching Traffic

Lookups move randomly in the network.

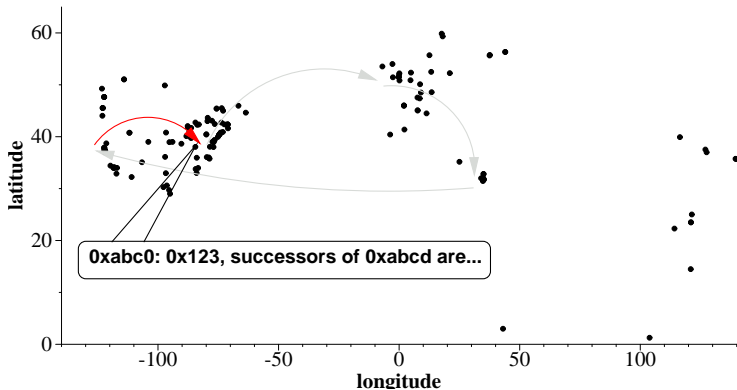
- Node IDs, block keys are assigned (pseudo-)randomly.



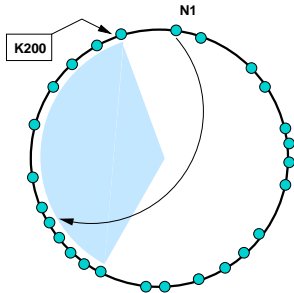
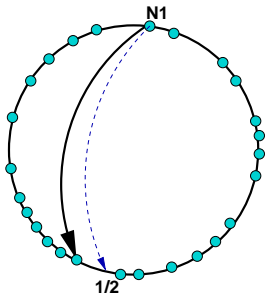
Watching Traffic

Lookups move randomly in the network.

- Node IDs, block keys are assigned (pseudo-)randomly.

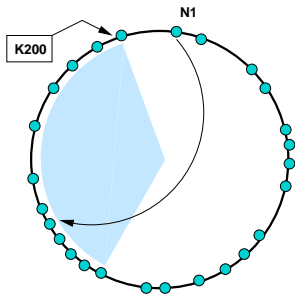


Chord: Shortening Lookup Hops Using Vivaldi



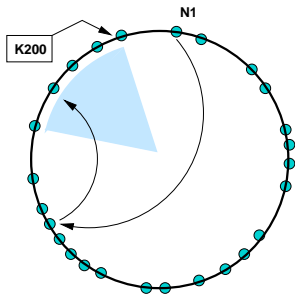
- Relax definition of finger from immediate successor to any node in range.
- Given more choices, can choose nearer neighbors.

Chord: Lookup Using Neighbor Selection



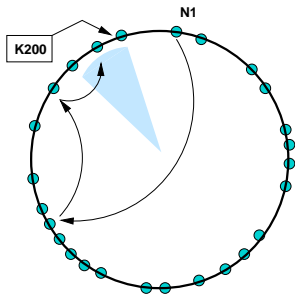
- Each hop now has a range of choices.

Chord: Lookup Using Neighbor Selection



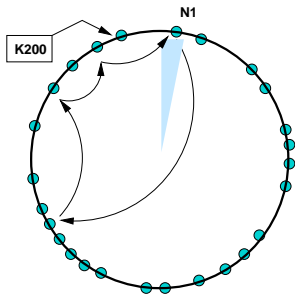
- Each hop now has a range of choices.

Chord: Lookup Using Neighbor Selection



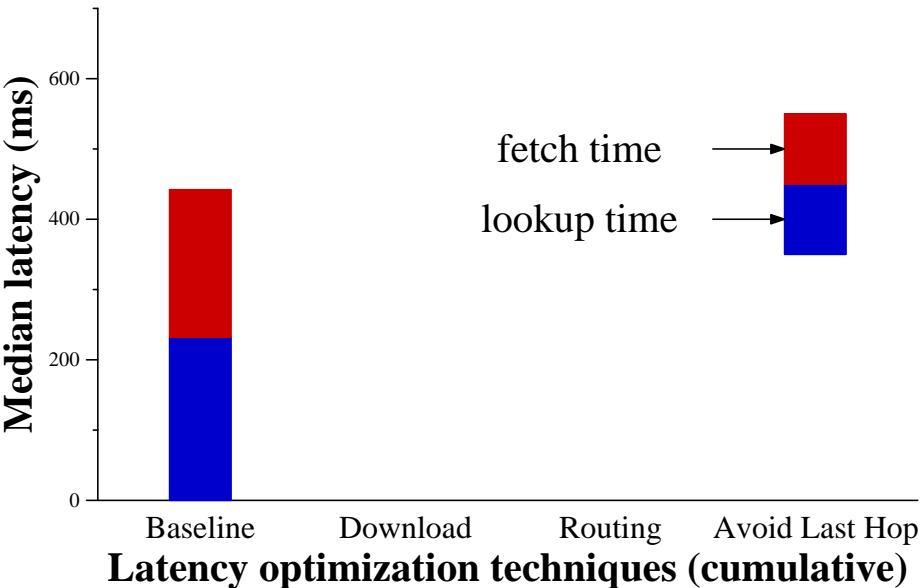
- Each hop now has a range of choices.

Chord: Lookup Using Neighbor Selection

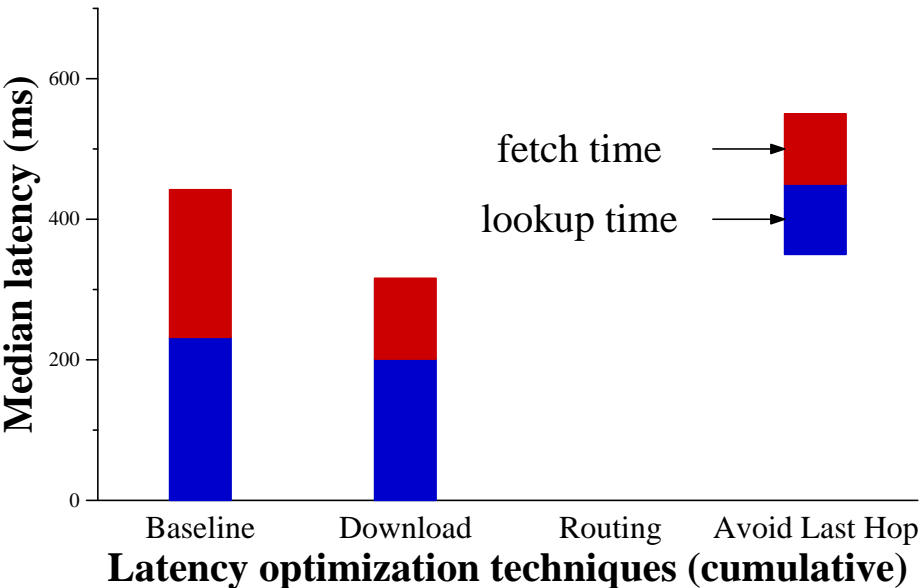


- Each hop now has a range of choices.

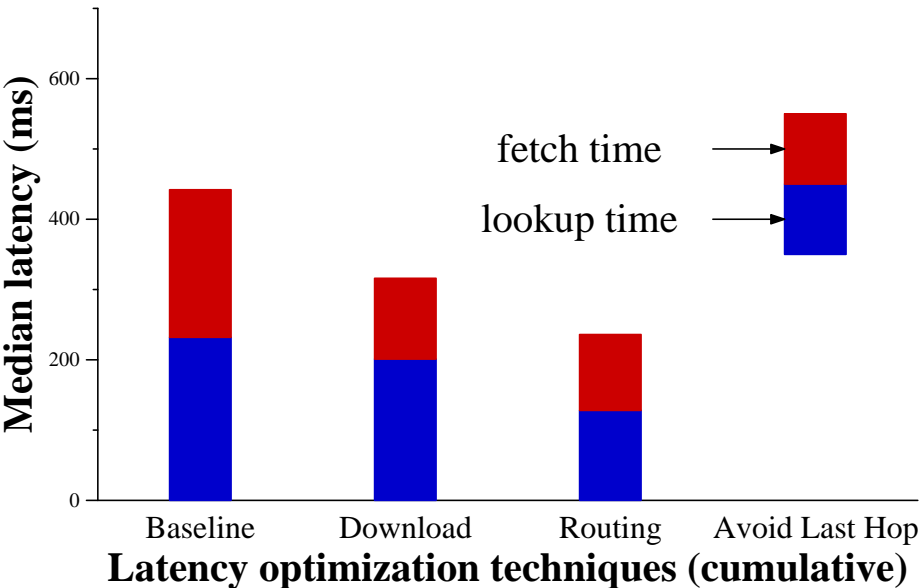
Vivaldi Improves Chord and DHash



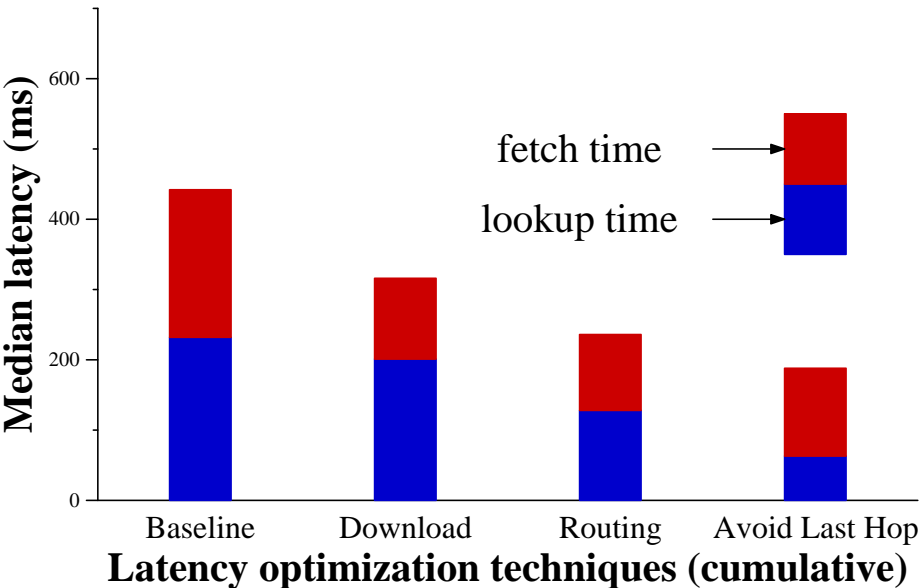
Vivaldi Improves Chord and DHash



Vivaldi Improves Chord and DHash



Vivaldi Improves Chord and DHash



Vivaldi Improves Chord and DHash

Vivaldi improved the performance of Chord and DHash:

- Chord lookups can use nearby neighbors instead of hopping all over the globe.
- DHash replicates blocks on r successors of block key.
- DHash can download block from nearest replica.

- Net effect: Vivaldi reduced block fetch time by close to 50% on PlanetLab.

Vivaldi is also used by other systems:

- Bamboo Distributed Hash Table
- SWORD Resource Discovery system
- others in development

- Other location techniques (IDMaps, IP2Geo) use static data (AS maps, guesses at physical location).
- Centralized coordinate systems (GNP, Lighthouse) need well-known landmark nodes.
- Decentralized coordinate systems (PIC, NPS) have been developed concurrently.
- Tang and Crovella (IMC 2003) analyze best Euclidean models to use; Shavitt and Tankel (Infocom 2004) suggest using hyperbolic geometries.

Chord is a scalable peer-to-peer system.

Vivaldi:

- accurately predicts round trip time *between node pairs not directly measured.*
- works without centralized infrastructure.
- improves the performance of a Chord and DHash

Height vectors are a promising coordinate space for the Internet.

Chord home page

- <http://pdos.lcs.mit.edu/chord>

Project IRIS (Peer-to-peer research)

- <http://project-iris.net>

Email

- rsc@mit.edu