

Proposal

Andrew Rosen

September 21, 2014

Contents

1	Hypotheses	3
1.1	DHTs are better for distributed computing under many circumstances . .	3
1.1.1	Really Cool Qualities	3
	Robustness and Fault-Tolerance	3
	Load Balancing	4
	Scalability	4
	Heterogeneity	5
	Ease of Adding, Removing, Maintaining, and money related factors	5
1.2	Different or subproblem: Certain DHTs are better at one application than another due to differences	5
1.2.1	Design Differences Impacts	5
1.2.2	Geometries	5
1.2.3	Routing Table Construction	5
1.2.4	Implementation Differences Impacts	5
	Recursive or iterative seek	5
1.2.5	MANETS	5
2	Justification and Why I Think It's Cool	6
2.0.6	DHTs well understood	6
2.0.7	DHTs are Highly used for their intended purposed	6
	Bittorrent, WoW	6
2.0.8	DHTs are being effectively leveraged for other things besides file sharing already	6
	PaaS	6
	Load Balancing in the cloud	6
	Computing is a natural extension	6
3	Possible Experiments and Applications	7
3.0.9	Map Reduce	7
	ChordReduce	7
	Comparison of MapReduce paradigm on different DHTs	7
3.0.10	High End Computing	7
	Metadata Management	7
	Robustness	7

3.0.11	Graph Processing on a DHT	7
	Embedding	7
	Distribute the work for solving a graph on a DHT	7
	Comparison to well established or state of the art methods . . .	7
3.0.12	Machine Learning Problems on A DHT	7
	Bayesian Learning	7
3.0.13	DHTs as a volunteer Platform	7
4	DHT Background	8

Chapter 1

Hypotheses

1.1 DHTs are better for distributed computing under many circumstances

Distributed Hash Tables (DHTs) are traditionally used as the backbone of P2P file-sharing applications.

However, they are seeing increasing use in other applications

- Citation(Hybrid computation)
-

1.1.1 Really Cool Qualities

Many DHTs are built upon the assumption that they will be used in some kind of P2P application.

Robustness and Fault-Tolerance

One of the most important assumptions of DHTs is that they are deployed on a non-static network. DHTs need to be built to account for a high level what is called *churn*. Churn refers to the disruption of routing caused by the constant joining and leaving of nodes.

This is mitigated by a few factors. First, no single node is a single point of failure for the network; this is accomplished by each node in the routing table having a small portion of the both the routing table and the information stored on the DHT (see the Load Balancing property below). Second is DHTs have a cheap repair or maintenance process that takes care of damage and integrating nodes that affects a minimal number of nodes, minimizing overhead. Third, DHTs often integrate a backup process into their protocols so that when a node goes down, one of the neighbors immediately take over its responsibility. The last property is that the Hashing Algorithm used to distribute content evenly across the network(again see load balancing) also distributes nodes evenly across the DHT. This means that nodes in the same geographic region occupy vastly different positions in the keyspace. If an entire geographic region is affected by a network outage,

Size	Speed	DHTs
$\log n$	$\log n$	Everybody

Table 1.1: The different ratios and their associated DHTs

this damage is spread evenly across the DHT, which can be manageable This property is the most important, as it deals with failure of entire sections of the network, rather than a single node.

So what are the consequences of these properties?

- DHTs are highly resilient to damage and often at higher rates of disruption that is normally.
-
-

Recent research in using DHTs for High End Computing [?] shows what can happen if we remove this assumption by placing the network that is almost completely static.

Load Balancing

All Distributed Hash Tables associate nodes and file identifiers with keys. These keys are generated by passing the identifiers into a hash function, typically SHA-160. The chosen hash function is typically large enough to avoid hash collisions and generates keys in a uniform manner. The result of this is that as more nodes join the network, the distribution of nodes in the keyspace becomes more uniform, as does the distribution of files.

While most DHTs follow this scheme, there is some minor variation on how keys are generated. Some DHTs can or do use geographic information to generate keys. In another instance, keys are static [?].

Scalability

How a DHT handles scalability is what causes the greatest variation between DHTs(HELP). Each DHT strikes a balance between the overhead cost of maintaining the routing table and the routing speed. The greatest architectural differences arise from the tradeoff of routing table size to lookup time and how it is maintained.

The vast majority of Distributed Hash Tables

Other factors: Bases, maintenense process, maintanense intervals,

Heterogeneity

Ease of Adding, Removing, Maintaining, and money related factors

1.2 Different or subproblem: Certain DHTs are better at one application than another due to differences

1.2.1 Design Differences Impacts

1.2.2 Geometries

1.2.3 Routing Table Construction

1.2.4 Implementation Differences Impacts

Recursive or iterative seek

1.2.5 MANETS

Chapter 2

Justification and Why I Think It's Cool

2.0.6 DHTs well understood

2.0.7 DHTs are Highly used for their intended purposed

Bittorrent, WoW

2.0.8 DHTs are being effectively leveraged for other things besides file sharing already

PaaS

Load Balancing in the cloud

Computing is a natural extension

Chapter 3

Possible Experiments and Applications

3.0.9 Map Reduce

ChordReduce

Comparison of MapReduce paradigm on different DHTs

3.0.10 High End Computing

Metadata Management

Robustness

3.0.11 Graph Processing on a DHT

Embedding

Distribute the work for solving a graph on a DHT

Comparison to well established or state of the art methods

3.0.12 Machine Learning Problems on A DHT

Bayesian Learning

3.0.13 DHTs as a volunteer Platform

Chapter 4

DHT Background

Bibliography