# Peer-to-Peer Prefix Tree for Large Scale Service Discovery

Cédric Tedeschi

October 2, 2008



DLPT : a Trie-Based Solution Mapping & Load Balancing October 2000 Summary & Open Problems

#### **Computing Needs**





#### Cosmology



Genomics

Nuclear security

#### **Computing Power**

## **Computational Grids**



#### Cédric Tedeschi

00			
0000	000000	000	



00			
0000	000000	000	



#### Cédric Tedeschi

0000	000000	0000	





0000	000000	0000	





0		



0		



000	000000	
0		



000	000000	
0		



000	000000	
0		



## P2P Information Retrieval (review)

#### Unstructured (Gnutella)

- Flooding
- Cost
- Non-exhaustiveness

#### Structured (DHTs)

- Routing
- Exhaustiveness
- Scalability
  - Logarithmic state
  - Logarithmic path
- Drawbacks
  - Exact queries only

## P2P Information Retrieval (review)

#### Unstructured (Gnutella)

- Flooding
- Cost
- Non-exhaustiveness

## Structured (DHTs)

- Routing
- Exhaustiveness
- Scalability
  - Logarithmic state
  - Logarithmic path
- Drawbacks
  - Exact queries only

## P2P Information Retrieval (review)

#### Unstructured (Gnutella)

- Flooding
- Cost
- Non-exhaustiveness

## Structured (DHTs)

- Routing
- Exhaustiveness
- Scalability
  - Logarithmic state
  - Logarithmic path
- Drawbacks
  - Exact queries only

000	000000	000	

#### Objectives

#### A P2P structured solution with

- Expressiveness
  - Range queries
  - Automatic completion
  - Multi-attribute queries
- Efficiency in a P2P environment
  - Heterogeneous ⇒ Load balancing
  - Dynamic ⇒ Fault-tolerance



# Outline

- DLPT : a Trie-Based Solution
- 2 Mapping & Load Balancing
  - Protocol
  - Load Balancing
  - Simulation
- 3

## Fault-tolerance

- Snap-Stabilizing PGCP Tree (State Model)
- Self-Stabilizing PGCP Tree (Message-Passing Model)
- Prototype Implementation
  - Design
  - Early Experiments
- 5 Summary & Open Problems

DLPT : a Trie-Based Solution			
	00 000 0		

# Outline

## DLPT : a Trie-Based Solution

- 2 Mapping & Load Balancing
  - Protocol
  - Load Balancing
  - Simulation
- 3 Fault-tolerance
  - Snap-Stabilizing PGCP Tree (State Model)
  - Self-Stabilizing PGCP Tree (Message-Passing Model)
- Prototype Implementation
  - Design
  - Early Experiments
- 5 Summary & Open Problems

ng & Load Balancing

Fault-tolerance Pr

Prototype Implementation Summary & Open Problems
Open

## DLPT : A trie-based indexing system

#### **Distributed Logical Structure**

Greatest Common Prefix Tree

Dynamically constructed

Bounded degree and height

#### Definition

Each node is the greatest common prefix of any pair of its children

oping & Load Balancing

Fault-tolerance Pr

Prototype Implementation Summary & Open Problems

## DLPT : A trie-based indexing system

#### Distributed Logical Structure

- Greatest Common Prefix Tree
- Dynamically constructed

Bounded degree and height



apping & Load Balancing o oo

Prototype Implementation Summary & Open Problems

## DLPT : A trie-based indexing system

#### Distributed Logical Structure

- Greatest Common Prefix Tree
- Dynamically constructed
- Bounded degree and height

DGEMM



& Load Balancing Fa

Fault-tolerance Pr 000000000 0 000000 00

Prototype Implementation Summary & Open Problems
Ooo

## DLPT : A trie-based indexing system

#### **Distributed Logical Structure**

Greatest Common Prefix Tree

Dynamically constructed

• Bounded degree and height

DGEMM

DGEMM DTRSM



pping & Load Balancing o

Fault-tolerance Pr

e Prototype Implementation Summary & Open Problems

## DLPT : A trie-based indexing system

#### **Distributed Logical Structure**

- Greatest Common Prefix Tree
- Dynamically constructed
- Bounded degree and height

DGEMM

O O DGEMM DTRSM



DLPT : a Trie-Based Solution Mapping & Load Balancing Geodesic Prototype Implementation Summary & Open Problem

## DLPT : A trie-based indexing system

#### **Distributed Logical Structure** Greatest Common Prefix Tree DTR DGEMM Dynamically constructed Bounded degree and height DGEMM DTRSM DGEMM DTRMM DTRSM Lookup clie Exact match S3L 1) Lookup (DTRMM) 5) val1 S3L mat mult val5 DGEMM 4 S3L trans

DTRMM DTRSM

S3L mat mult noadd

DLPT : a Trie-Based Solution Mapping & Load Balancing Fault-tolerance Prototype Implementation Summary & Open Problem

## DLPT : A trie-based indexing system



#### Cédric Tedeschi

DLPT : a Trie-Based Solution Mapping & Load Balancing Fault-tolerance Prototype Implementation Summary & Open Problem

## DLPT : A trie-based indexing system



#### Cédric Tedeschi



# Outline

## DLPT : a Trie-Based Solution

- 2 Mapping & Load Balancing
  - Protocol
  - Load Balancing
  - Simulation
- 3) Fault-tolerance
  - Snap-Stabilizing PGCP Tree (State Model)
  - Self-Stabilizing PGCP Tree (Message-Passing Model)
- Prototype Implementation
  - Design
  - Early Experiments
- 5 Summary & Open Problems



## The Mapping Problem







Cédric Tedeschi

A Peer-to-Peer Prefix Tree for Large Scale Service Discovery 10/44



## The Mapping Problem

- Associate a **Node** to a **Peer**
- An underlying DHT



DLPT : a Trie-Based Solution Mapping & Load Balancing Fault-tolerance Prototype Implementation Summary & Open Problems

# The Mapping Problem

## Drawbacks

Need for a DHT

- Costly
- Random mapping (No clustering)



10/44

Mapping & Load Balancing		
•0		
0		



- Build a ring over peers
- Each peer maintains at least one node
- Map using consistent hashing
- Clustering

Mapping & Load Balancing		
•0		
0		



DTRMM DTRSM

#### Build a ring over peers

- Each peer maintains at least one node
- Map using consistent hashing
- Clustering

Mapping & Load Balancing		
•0		
0		



DTRMM DTRSM

- Build a ring over peers
- · Each peer maintains at least one node
- Map using consistent hashing
- Clustering



DGEMM

PROD DTRMM Mat DTRSM

DTRMM DTRSM

DGEMM

D

DTR

Build a ring over peers

Mat

MatPROD

MatSUM

Each peer maintains at least one node

RsBFz

MatSUM

MatPROD

Lrces

- Map using consistent hashing
- Clustering

DMr2

DTRMM DTRSM

Lrces

MatSUM

MatPROD

DGEMM

RsBFz

DMrXz

Mapping & Load Balancing ○● ○○○

Fault-tolerance I

e Prototype Implementation Summary & Open Problems

## Insertion of a New Peer



# Join Algorithm Random contact peer Initiating the tree routing Routing in the tree Successor of *P* is *P<sub>n</sub>* or *succ*(*P<sub>n</sub>*)

Mapping & Load Balancing ○● ○○

Fault-tolerance I

e Prototype Implementation Summary & Open Problems

## Insertion of a New Peer





Mapping & Load Balancing ○● ○○

Fault-tolerance

ce Prototype Implementation Summary & Open Problems

## Insertion of a New Peer



Join	Algorithm
1	Random contact peer
2	Initiating the tree routing
3	Routing in the tree
	Successor of $P$ is $P_{p}$ or $succ(P_{p})$
Fault-tolerance F

ce Prototype Implementation Summary & Open Problems

### Insertion of a New Peer



Join Algorithm	
<ol> <li>Random contact peer</li> </ol>	
Initiating the tree routing	
3 Routing in the tree	

4 Successor of P is  $P_n$  or  $succ(P_n)$ 

### Insertion of a New Peer



00

#### Join Algorithm

- Random contact peer
- Initiating the tree routing 2
- Routing in the tree (3
- Successor of *P* is  $P_n$  or  $succ(P_n)$

### Join Complexities

#### O(Trie complexities)

# Insertion of a New Peer



000

#### Join Algorithm

- Random contact peer
- Initiating the tree routing
- Routing in the tree (3
- Successor of *P* is  $P_n$  or  $succ(P_n)$

### Join Complexities

O(Trie complexities)

### Load balancing

- Depth of nodes
- Popularity of services
- Capacity of peers

# Load balancing heuristics - related work (DHTs)

Each peer is assigned a set of items

- Karger and Ruhl, 2001
  - · Periodic random item balancing
  - Homogeneity of peer capacities
- Godfrey et al., 2003
  - Periodic item redistribution
  - Semi-centralized
- Ledlie and Seltzer, 2005
  - Chooses the best location for a joining peer among k
  - Heterogeneity of peer capacities and data popularity

Fault-toleranceP000000000000000000

Prototype Implementation Summary & Open Problems

# A novel heuristic : Max Local Throughput



End of period  $\tau$ S get load information from P

Fault-tolerance F

Prototype Implementation Summary & Open Problems



End of period $ au$	Search for a new	
S get load information	locally optimal	
from P	distribution	

Fault-toleranceP000000000000000000

e Prototype Implementation Summary & Open Problems





- At the end of period  $\tau$
- Two adjacent peers S and P with capacity  $C_S$  and  $C_P$
- $\nu_S^{\tau}$  and  $\nu_P^{\tau}$  the sets of nodes currently managed by S and P
- $L_S^{\tau} = \sum_{n \in \nu_S^{\tau}} I_n$
- $L_P^{\tau} = \sum_{n \in \nu_P^{\tau}} I_n$
- $T_{S,P}^{\tau} = min(L_S^{\tau}, C_S) + min(L_P^{\tau}, C_P)$
- Maximize the throughput of au+1 based on information of au
- Find  $\nu_S^{\tau+1}$  and  $\nu_P^{\tau+1}$  that maximizes  $T_{SP}^{\tau+1}$ .

- At the end of period  $\tau$
- Two adjacent peers S and P with capacity  $C_S$  and  $C_P$
- $\nu_{S}^{\tau}$  and  $\nu_{P}^{\tau}$  the sets of nodes currently managed by S and P
- $L_{\mathcal{S}}^{\tau} = \sum_{n \in \nu_{\mathcal{S}}^{\tau}} I_n$
- $L_P^{\tau} = \sum_{n \in \nu_P^{\tau}} I_n$
- $T_{S,P}^{\tau} = min(L_S^{\tau}, C_S) + min(L_P^{\tau}, C_P)$
- Maximize the throughput of au+1 based on information of au
- Find  $\nu_S^{\tau+1}$  and  $\nu_P^{\tau+1}$  that maximizes  $T_{S,P}^{\tau+1}$ .

- At the end of period  $\tau$
- Two adjacent peers S and P with capacity  $C_S$  and  $C_P$
- $\nu_{S}^{\tau}$  and  $\nu_{P}^{\tau}$  the sets of nodes currently managed by S and P
- $L_{\mathcal{S}}^{\tau} = \sum_{n \in \nu_{\mathcal{S}}^{\tau}} I_n$
- $L_P^{\tau} = \sum_{n \in \nu_P^{\tau}} I_n$
- $T^{\tau}_{\mathcal{S},\mathcal{P}} = \min(L^{\tau}_{\mathcal{S}},\mathcal{C}_{\mathcal{S}}) + \min(L^{\tau}_{\mathcal{P}},\mathcal{C}_{\mathcal{P}})$
- Maximize the throughput of au + 1 based on information of au
- Find  $\nu_S^{\tau+1}$  and  $\nu_P^{\tau+1}$  that maximizes  $T_{S,P}^{\tau+1}$ .

- At the end of period  $\tau$
- Two adjacent peers S and P with capacity  $C_S$  and  $C_P$
- $\nu_{S}^{\tau}$  and  $\nu_{P}^{\tau}$  the sets of nodes currently managed by S and P
- $L_{\mathcal{S}}^{\tau} = \sum_{n \in \nu_{\mathcal{S}}^{\tau}} I_n$
- $L_P^{\tau} = \sum_{n \in \nu_P^{\tau}} I_n$
- $T^{\tau}_{\mathcal{S},\mathcal{P}} = \min(L^{\tau}_{\mathcal{S}},\mathcal{C}_{\mathcal{S}}) + \min(L^{\tau}_{\mathcal{P}},\mathcal{C}_{\mathcal{P}})$
- Maximize the throughput of au+1 based on information of au
- Find  $\nu_S^{\tau+1}$  and  $\nu_P^{\tau+1}$  that maximizes  $T_{S,P}^{\tau+1}$ .

- At the end of period  $\tau$
- Two adjacent peers S and P with capacity  $C_S$  and  $C_P$
- $\nu_{S}^{\tau}$  and  $\nu_{P}^{\tau}$  the sets of nodes currently managed by S and P
- $L_{\mathcal{S}}^{\tau} = \sum_{n \in \nu_{\mathcal{S}}^{\tau}} I_n$
- $L_P^{\tau} = \sum_{n \in \nu_P^{\tau}} I_n$
- $T^{\tau}_{\mathcal{S},\mathcal{P}} = min(L^{\tau}_{\mathcal{S}},\mathcal{C}_{\mathcal{S}}) + min(L^{\tau}_{\mathcal{P}},\mathcal{C}_{\mathcal{P}})$
- Maximize the throughput of au+1 based on information of au
- Find  $\nu_{S}^{\tau+1}$  and  $\nu_{P}^{\tau+1}$  that maximizes  $T_{S,P}^{\tau+1}$ .

# DLPT : a Trie-Based Solution Mapping & Load Balancing Fault-tolerance Prototype Implementation Summary & Open Problem

### Simulation results

Load	Stable network		Dynamic network	
	Max local th.	K-choices	Max local th.	K-choices
5%	39,62%	38,58%	18.25%	32,47%
10%	103,41%	58,95%	46,16%	51,00%
16%	147,07%	64,97%	65,90%	59,11%
24%	165,25%	59,27%	71,26%	60,01%
40%	206,90%	68,16%	97,71%	67,18%
80%	230,51%	76,99%	90,59%	71,93%

# 

### Simulation results

Load balancing - dynamic network - dynamic load Max Local Throughput [50 run] K-choices [50 run No LB 50 run 0.8 Satisfied requests 0.6 0.4 0.2 0 0 20 40 60 80 100 120 140 160 Time



### Simulation results



	Fault-tolerance		
0000	000000	000	

# Outline

- DLPT : a Trie-Based Solution
- 2 Mapping & Load Balancing
  - Protocol
  - Load Balancing
  - Simulation
- 3 Fault-tolerance
  - Snap-Stabilizing PGCP Tree (State Model)
  - Self-Stabilizing PGCP Tree (Message-Passing Model)
- Prototype Implementation
  - Design
  - Early Experiments
  - Summary & Open Problems

	Fault-tolerance	
00 000 0		

### • Creation and maintenance of the tree

### Fault-tolerance

- Usually, based on replication
  - Costly
  - Unable to recover after arbitrary transient failures
- Best-effort ?

	Fault-tolerance	
00 000 0		

- Creation and maintenance of the tree
- Fault-tolerance
  - Usually, based on replication
    - Costly
    - Unable to recover after arbitrary transient failures
  - Best-effort?

	Fault-tolerance	
00 000 0		

- Creation and maintenance of the tree
- Fault-tolerance
  - Usually, based on replication
    - Costly
    - Unable to recover after arbitrary transient failures
  - Best-effort?

	Fault-tolerance	
0		

### Self-Stabilization (Dijkstra, 74)

- General technique to tolerate transient faults
- Guaranteed to converge to the intended behavior
  - · Regardless of the initial state of the system
  - In finite time



DLPT : a Trie-Based Solution Mapping & Load Balancing Fault-tolerance Prototype Implementation Summary & Open Problems

### Architecture Model

### **Physical Network**

### Basic entity : peer (processor)

- *P*<sub>1</sub> can communicate with *P*<sub>2</sub> if *P*<sub>1</sub> *knows P*<sub>2</sub>.
- Runs some logical nodes

### Logical Tree

Basic entity : (tree) node

- Distributed among peers
- Mapping details abstracted out
- Protocol run inside nodes





- Message exchanges modeled by the ability to read variables of other nodes (**Neighbors**)
- A node can only write its own variables
- The state of a node is defined by the values of its variables
- The configuration of the system is the product of the states
- Actions : < guard > -> < statement >
  - < guard > bool. expr. of variables of p and its neighbors
  - statement >
    - Executed only if its < guard >= true
    - Updates one or more variables of p
- Distributed and unfair scheduler daemon



- Message exchanges modeled by the ability to read variables of other nodes (**Neighbors**)
- A node can only write its own variables
- The state of a node is defined by the values of its variables
- The configuration of the system is the product of the states
- Actions : < guard >→< statement >
  - < guard > bool. expr. of variables of p and its neighbors
  - < statement >
    - Executed only if its < guard >= true
    - Updates one or more variables of p
- Distributed and unfair scheduler daemon



- Message exchanges modeled by the ability to read variables of other nodes (**Neighbors**)
- A node can only write its own variables
- The state of a node is defined by the values of its variables
- The configuration of the system is the product of the states
- Actions : < guard >→< statement >
  - < guard > bool. expr. of variables of p and its neighbors
  - < statement >
    - Executed only if its < guard >= true
    - Updates one or more variables of p
- Distributed and unfair scheduler daemon



- Message exchanges modeled by the ability to read variables of other nodes (**Neighbors**)
- A node can only write its own variables
- The state of a node is defined by the values of its variables
- The configuration of the system is the product of the states
- Actions : < guard >→< statement >
  - < guard > bool. expr. of variables of p and its neighbors
  - < statement >
    - Executed only if its < guard >= true
    - Updates one or more variables of p
- Distributed and unfair scheduler daemon

	Fault-tolerance	
00 000 0	00000000 000000	

### Snap-stabilization

Let P be a protocol designed to solve a task T. P is said **snap-stabilizing** if and only if, starting from any configuration, any execution e of P always satisfies the specification of T.

### The Distributed Structures Maintained





### **PGCP** Tree

A labeled rooted tree s.t. each node label is the **PGCP** of any pair of its children labels.

#### Prefix Heap

A labeled rooted tree s.t. each node label is the **PGCP** of all its children labels.

### The Distributed Structures Maintained



A Prefix Heap s.t. for each node p, every pair of children (q, r):

$$\begin{cases} (1) \quad l_q \neq l_r \\ (2) \quad l_q(\text{resp. } l_r) \text{ is not a prefix of } l_r(\text{resp. } l_q) \\ (3) \quad |GCP(l_q, l_r)| = |l_p| \end{cases}$$

is a PGCP Tree.

# Snap-stabilizing PGCP Tree construction

### **Topology Assumptions**

The initial graph is a rooted connected tree.

### Snap-stabilizing PGCP Tree

A protocol P is considered as a snap-stabilizing PGCP tree construction if and only if, any execution initiated by the root terminates in finite time and when the execution terminated, the tree is a PGCP tree.

DLPT : a Trie-Based Solution Mapping & Load Balancing Solution October October

# The Protocol

### **Functions**

- newNode(lbl, st, chldn)
- destroy (p)
- heapify() locally creates a heap
- repair() locally builds a PGCP Tree from a heap

#### Phases

- From the snap-stabilizing PIF [Bui, Datta, Petit, Villain 99]
- Broadcast phase : top-down wave initiating the algorithm
- Heapify phase : down-top traversal building a prefix heap
- Repair phase : final top-down wave building a PGCP tree

DLPT : a Trie-Based Solution Mapping & Load Balancing Sault-tolerance October October

# The Protocol

#### Functions

- newNode(lbl, st, chldn)
- destroy (p)
- heapify() locally creates a heap
- repair() locally builds a PGCP Tree from a heap

#### Phases

- From the snap-stabilizing PIF [Bui, Datta, Petit, Villain 99]
- Broadcast phase : top-down wave initiating the algorithm
- Heapify phase : down-top traversal building a prefix heap
- Repair phase : final top-down wave building a PGCP tree


















































	Fault-tolerance	
	000000000	
0		



















	Fault-tolerance	
00	000000000 000000	



	Fault-tolerance	
	000000000 000000	
0		



	Fault-tolerance	
00 000 0	000000000 000000	



	Fault-tolerance	
	000000000 000000	
0		



	Fault-tolerance	
00	000000000 000000	
0		



	Fault-tolerance	
	000000000 000000	
0		



	Fault-tolerance	
	000000000 000000	
0		



Mapping & Load Balancing	Fault-tolerance	

### Complexities

- Average convergence time to PGCP tree : O(h+h') rounds
- Worst case :
  - O(n) rounds
  - O(n<sup>2</sup>) operations
  - O(n) extra space (degree)

# Simulation : Time Complexity



27/44

## Simulation : Extra Space



28/44

# Snap in State Model vs Self in Message Passing

#### Snap : Drawbacks

- Assumes a rooted connected tree
- Written in a coarse grain communication model

#### A New Protocol

- Work with any initial topology
- Designed in the message passing model

# Snap in State Model vs Self in Message Passing

#### Snap : Drawbacks

- Assumes a rooted connected tree
- Written in a coarse grain communication model

#### A New Protocol

- Work with any initial topology
- Designed in the message passing model

	Fault-tolerance	
00 000 0	000000000000000000000000000000000000000	

### Steps

- · Launched periodically on each node
- 1 Parent processing
- 2 Children processing

	Fault-tolerance	
00 000 0	000000000	



	Fault-tolerance	
00 000 0	000000000	



	Fault-tolerance	
00 000 0	000000000000000000000000000000000000000	



	Fault-tolerance	
00 000 0	000000000000000000000000000000000000000	



	Fault-tolerance	
00 000 0	000000000	



	Fault-tolerance	
00 000 0	000000000	



	Fault-tolerance	
00 000 0	000000000000000000000000000000000000000	



	Fault-tolerance	
00 000 0	000000000000000000000000000000000000000	



	Fault-tolerance	
00 000 0	000000000000000000000000000000000000000	



	Fault-tolerance	
00 000 0	000000000	


	Fault-tolerance	
00 000 0	000000000	



	Fault-tolerance	
00 000 0	000000000	



	Fault-tolerance	
00 000 0	000000000000000000000000000000000000000	



	Fault-tolerance	
00 000 0	000000000000000000000000000000000000000	



	Fault-tolerance	
00 000 0	000000000000000000000000000000000000000	



	Fault-tolerance	
00 000 0	000000000000000000000000000000000000000	



	Fault-tolerance	
00 000 0	000000000	



	Fault-tolerance	
00 000 0	000000000000000000000000000000000000000	



	Fault-tolerance	
00 000 0	000000000	



	Fault-tolerance	
00 000 0	000000000000000000000000000000000000000	



# Simulation : Convergence Time



Cédric Tedeschi

A Peer-to-Peer Prefix Tree for Large Scale Service Discovery

DLPT : a Trie-Based Solution Mapping & Load Balancing Pault-tolerance Ocoobooo Ocooboo O

### Simulation : Communication Amount



Cédric Tedeschi

A Peer-to-Peer Prefix Tree for Large Scale Service Discovery

	Fault-tolerance	
	000000	
0		

### Simulation : Degradation



Cédric Tedeschi

	Prototype Implementation	
0		

# Outline

- DLPT : a Trie-Based Solution
- 2 Mapping & Load Balancing
  - Protocol
  - Load Balancing
  - Simulation
- 3 Fault-tolerance
  - Snap-Stabilizing PGCP Tree (State Model)
  - Self-Stabilizing PGCP Tree (Message-Passing Model)
- Prototype Implementation
  - Design
  - Early Experiments
  - Summary & Open Problems

	Prototype Implementation	
00 000 0	• 000	

#### Physical Layer

- JXTA-based
- JXTA discovery
- "DLPT" JXTA-group



# DLPT : a Trie-Based Solution Mapping & Load Balancing Control Control

#### Logical Layer

- Explicit tree topology maintenance
- Random distribution of nodes on peers
- Message in the tree :
  - Serialized communication between JXTA peers
  - JXTA communications

#### **Physical Layer**

- JXTA-based
- JXTA discovery
- "DLPT" JXTA-group



# **Experiments : Set-Up**

#### Platform used

#### Grid'5000 platform

- Capricorne cluster in Lyon
  - 25 nodes (2 x AMD Opteron 246)
  - Gigabit Ethernet (Myrinet-2000 cards)
- Grelon cluster in Nancy
  - 60 nodes (2 x Intel Xeon 5110)
  - Gigabit Ethernet (Broadcom BCM5721 cards)

#### Steps

- Physical deployment
- Simulation of servers (insertion requests)
- Simulation of clients (discovery requests)



### Experiments : Capricorne Results



38/44

DLPT : a Trie-Based Solution Mapping & Load Balancing Goo Ocooco Ocooco

### Experiments : Grelon Results

Nodes (Services)	Peers	Response time (s)
432 (300)	60	0.105

		Summary & Open Problems
0		

# Outline

- DLPT : a Trie-Based Solution
- 2 Mapping & Load Balancing
  - Protocol
  - Load Balancing
  - Simulation
- 3 Fault-tolerance
  - Snap-Stabilizing PGCP Tree (State Model)
  - Self-Stabilizing PGCP Tree (Message-Passing Model)
- Prototype Implementation
  - Design
  - Early Experiments

### 5 Summary & Open Problems



- A novel P2P-fashioned approach for service discovery
- Structured and Expressive
- Mapping and load balancing
  - Self-contained architecture
  - Load balancing heuristics
- Fault-tolerance
  - Best-effort alternatives
  - Snap-stabilizing protocol in the state model
  - Self-stabilizing protocols in message-passing
- Software Prototype
  - Two layer architecture implemented
  - Promising first experiments



- A novel P2P-fashioned approach for service discovery
- Structured and Expressive
- Mapping and load balancing
  - Self-contained architecture
  - Load balancing heuristics
- Fault-tolerance
  - Best-effort alternatives
  - Snap-stabilizing protocol in the state model
  - Self-stabilizing protocols in message-passing
- Software Prototype
  - Two layer architecture implemented
  - Promising first experiments



- A novel P2P-fashioned approach for service discovery
- Structured and Expressive
- Mapping and load balancing
  - Self-contained architecture
  - Load balancing heuristics
- Fault-tolerance
  - Best-effort alternatives
  - Snap-stabilizing protocol in the state model
  - Self-stabilizing protocols in message-passing
- Software Prototype
  - Two layer architecture implemented
  - Promising first experiments



- A novel P2P-fashioned approach for service discovery
- Structured and Expressive
- Mapping and load balancing
  - Self-contained architecture
  - Load balancing heuristics
- Fault-tolerance
  - Best-effort alternatives
  - Snap-stabilizing protocol in the state model
  - Self-stabilizing protocols in message-passing
- Software Prototype
  - Two layer architecture implemented
  - Promising first experiments

- Design
  - Topology awareness
  - Optimizations (cache, shortcuts)
- Mapping & Load balancing
  - Coupling our objective function with K-choices
  - Analysis results on the clustering obtained
  - Implementation
- Fault-tolerance
  - Snap-stabilization and message-passing
  - Self-stabilization and availability
  - Implementation
- Experiments
  - Load balancing
  - Self-stabilization
  - Larger scale

DLPT : a Trie-Based Solution Mapping & Load Balancing Generation Summary & Open Problems

- Design
  - Topology awareness
  - Optimizations (cache, shortcuts)
- Mapping & Load balancing
  - Coupling our objective function with K-choices
  - Analysis results on the clustering obtained
  - Implementation
- Fault-tolerance
  - Snap-stabilization and message-passing
  - Self-stabilization and availability
  - Implementation
- Experiments
  - Load balancing
  - Self-stabilization
  - Larger scale

DLPT : a Trie-Based Solution Mapping & Load Balancing Generation Summary & Open Problems

- Design
  - Topology awareness
  - Optimizations (cache, shortcuts)
- Mapping & Load balancing
  - Coupling our objective function with K-choices
  - · Analysis results on the clustering obtained
  - Implementation
- Fault-tolerance
  - Snap-stabilization and message-passing
  - Self-stabilization and availability
  - Implementation
- Experiments
  - Load balancing
  - Self-stabilization
  - Larger scale

DLPT : a Trie-Based Solution Mapping & Load Balancing Generation Summary & Open Problems

- Design
  - Topology awareness
  - Optimizations (cache, shortcuts)
- Mapping & Load balancing
  - Coupling our objective function with K-choices
  - Analysis results on the clustering obtained
  - Implementation
- Fault-tolerance
  - Snap-stabilization and message-passing
  - Self-stabilization and availability
  - Implementation
- Experiments
  - Load balancing
  - Self-stabilization
  - Larger scale

 DLPT : a Trie-Based Solution
 Mapping & Load Balancing
 Fault-tolerance
 Prototype Implementation
 Summary & Open Problems

# Service Discovery : Open Problems

#### Different kinds of platforms

- Petascale computing
- Research grids
- Desktop computing

#### Different kinds of volatility

- Failures
- Crashes
- Resources sharing

		Summary & Open Problems
00 000 0		

# Thx !

Cédric Tedeschi

A Peer-to-Peer Prefix Tree for Large Scale Service Discovery 44/44