



Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

Sequencing

Transformation

Experiments

Conclusion

Efficient Similarity Search for Tree-Structured Data

Guoliang Li Xuhui Liu Jianhua Feng Lizhu Zhou
(SSDBM 2008)

Department of Computer Science and Technology, Tsinghua
University, Beijing 100084, China

July 10, 2008



RoadMap

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

Sequencing

Transformation

Experiments

Conclusion

- 1 Motivation
 - Problem
 - Applications
 - Similarity Metrics
 - Existing Methods
- 2 Sequence-Based Tree-Structured Data Similarity Search
 - Intuition
 - Sequencing
 - Edit Distance Transformation
- 3 Experiments
- 4 Conclusion



RoadMap

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

Sequencing

Transformation

Experiments

Conclusion

1 Motivation

- Problem
- Applications
- Similarity Metrics
- Existing Methods

2 Sequence-Based Tree-Structured Data Similarity Search

- Intuition
- Sequencing
- Edit Distance Transformation

3 Experiments

4 Conclusion



Problem

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

Sequencing

Transformation

Experiments

Conclusion

Problem

- Given a set of trees and a query tree
- Find all the trees that are ***similar*** to the query tree



Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

Sequencing

Transformation

Experiments

Conclusion

Applications

- Comparison of hierarchically structured data
- Alignment of RNA secondary structures in computational biology
- Approximate XML document match
- Schema mapping of tree-structured data



Tree Similarity Metrics

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

Sequencing

Transformation

Experiments

Conclusion

Tree Similarity Metrics

- Largest Common Subtree
- Smallest Common Super-tree
- Tree Edit Distance
- ...



Largest Common Subtree Distance

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

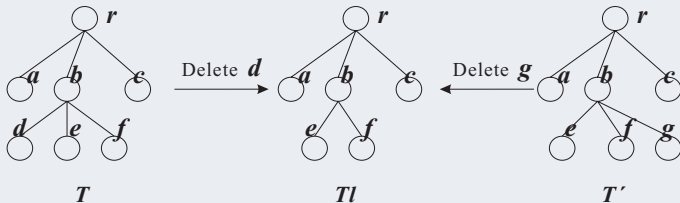
Sequencing

Transformation

Experiments

Conclusion

Largest Common Subtree



Largest Common Subtree Distance (LCST)

- the sum of # of operations to transfer the two trees into the largest common subtree



Smallest Common Super-tree Distance

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

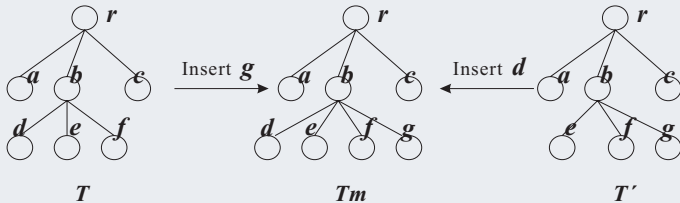
Sequencing

Transformation

Experiments

Conclusion

Smallest Common Super-tree



Smallest Common Super-tree Distance (SCST)

- the sum of # of operations to transfer the two trees into the smallest common super-tree



Tree Edit Distance

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

Sequencing

Transformation

Experiments

Conclusion

Tree Edit Operations

- Insert Node
- Delete Node
- Substitute Node

Tree Edit Distance

- # of tree edit operations



Edit Operations

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

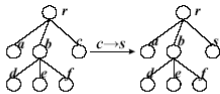
Sequencing

Transformation

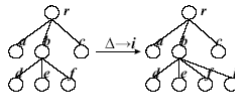
Experiments

Conclusion

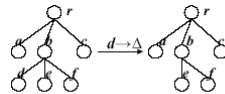
Examples



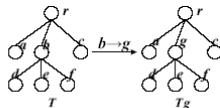
(a) Substitution of Leaf Node



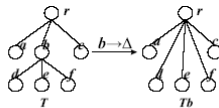
(b) Insertion of Leaf Node



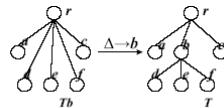
(c) Deletion of Leaf Node



(d) Substitution of Internal Node



(e) Deletion of Internal Node



(f) Insertion of Internal Node



Edit Operations

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

Sequencing

Transformation

Experiments

Conclusion

Edit Operations

$$\lambda_s(v) = \text{cSize}(v)$$

$$\lambda_d(v) = \begin{cases} 2 & \text{if } v \text{ is a leaf node} \\ 1 & \text{if } v \text{ is an internal node and } \text{parent}(v) = v \\ \text{cSize}(v) & \text{if } v \text{ is an internal node and } \text{parent}(v) \neq v \end{cases} \quad (1)$$

$$\lambda_i(v) = \begin{cases} 2 & \text{if } v \text{ is a leaf node} \\ 1 & \text{if } v \text{ is an internal node and } \text{parent}(v) = v \\ \text{cSize}(v) & \text{if } v \text{ is an internal node and } \text{parent}(v) \neq v \end{cases}$$



Tree Edit Distance

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

Sequencing

Transformation

Experiments

Conclusion



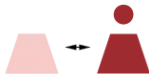
Substitution of l into l'

Deletion of l

Insertion of l'



$$\begin{aligned} &\text{Distance}(f, f') \\ &+ \\ &\text{sub}(l, l') \end{aligned}$$



$$\begin{aligned} &\text{Distance}(f, l'(f')) \\ &+ \\ &\text{del}(l) \end{aligned}$$



$$\begin{aligned} &\text{Distance}(l(f), f') \\ &+ \\ &\text{ins}(l') \end{aligned}$$



Forest Edit Distance

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

Sequencing

Transformation

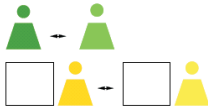
Experiments

Conclusion

$$l(f) \circ t \quad \begin{array}{c} \bullet \\ \text{Green Triangle} \end{array} \square \begin{array}{c} \bullet \\ \text{Yellow Triangle} \end{array} \rightarrow \begin{array}{c} \bullet \\ \text{Green Triangle} \end{array} \square \begin{array}{c} \bullet \\ \text{Yellow Triangle} \end{array} \quad l'(f') \circ t'$$

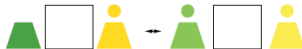
Leftmost decomposition

Substitution of l into l'



$$\begin{aligned} &\text{Distance}(l(f), l'(f')) \\ &\quad + \\ &\text{Distance}(t, t') \end{aligned}$$

Deletion of l



$$\text{Distance}(f \circ t, l'(f') \circ t') + \text{del}(l)$$

Insertion of l'



$$\text{Distance}(l(f) \circ t, f' \circ t') + \text{ins}(l')$$



Forest Edit Distance

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

Sequencing

Transformation

Experiments

Conclusion

$$t \circ l(f) \quad \begin{array}{c} \text{green} \\ \text{triangle} \end{array} \square \begin{array}{c} \text{yellow} \\ \text{triangle} \end{array} \leftrightarrow \begin{array}{c} \text{green} \\ \text{triangle} \end{array} \square \begin{array}{c} \text{yellow} \\ \text{triangle} \end{array} \quad t' \circ l'(f')$$

Rightmost decomposition

Substitution of l into l'



$$\text{Distance}(l(f), l'(f')) \\ + \\ \text{Distance}(t, t')$$



Deletion of l



$$\text{Distance}(t \circ f, t' \circ l'(f')) + \text{del}(l)$$

Insertion of l'



$$\text{Distance}(t \circ l(f), t' \circ f') + \text{ins}(l')$$



Forest Edit Distance

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

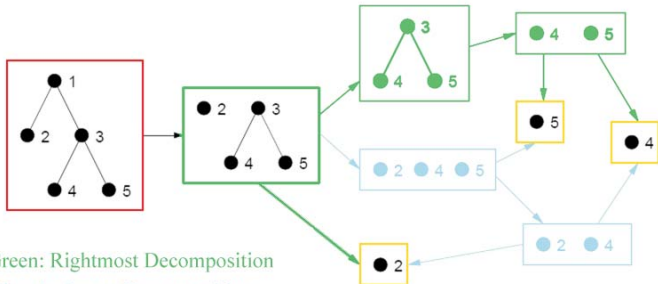
Intuition

Sequencing

Transformation

Experiments

Conclusion



Green: Rightmost Decomposition

Blue: Leftmost Decomposition



Tree Edit Distance

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

Sequencing

Transformation

Experiments

Conclusion

Edit Distance

- Different decomposition strategies
- Dynamic programming
- The costs of a commonly used algorithm
 - Space: $|T_1| * |T_2|$
 - Time:
 $|T_1| * |T_2| * \min(|depth(T_1)|, |leaves(T_1)|) * \min(|depth(T_2)|, |leaves(T_2)|)$
 - Worst Case: $|T_1|^2 * |T_2|^2$
- High CPU and IO costs!



Tree Edit Distance

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

Sequencing

Transformation

Experiments

Conclusion

Edit Distance

- Different decomposition strategies
- Dynamic programming
- The costs of a commonly used algorithm
 - Space: $|T_1| * |T_2|$
 - Time:
 $|T_1| * |T_2| * \min(|depth(T_1)|, |leaves(T_1)|) * \min(|depth(T_2)|, |leaves(T_2)|)$
 - Worst Case: $|T_1|^2 * |T_2|^2$
- High CPU and IO costs!



Complexity of Tree Similarity Metrics

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

Sequencing

Transformation

Experiments

Conclusion

Complexity

Tree edit distance

variant	type	time	space	reference
general	O	$O(T_1 T_2 D_1^2D_2^2)$	$O(T_1 T_2 D_1^2D_2^2)$	[43]
general	O	$O(T_1 T_2 \min(L_1, D_1)\min(L_2, D_2))$	$O(T_1 T_2)$	[55]
general	O	$O(T_1 ^2 T_2 \log T_2)$	$O(T_1 T_2)$	[25]
general	O	$O(T_1 T_2 + L_1^2 T_2 + L_1^2L_2)$	$O((T_1 + L_1^2)\min(L_2, D_2) + T_2)$	[8]
general	U	MAX SNP-hard		[54]
constrained	O	$O(T_1 T_2)$	$O(T_1 T_2)$	[51]
constrained	O	$O(T_1 T_2 I_1I_2)$	$O(T_1 D_2I_2)$	[37]
constrained	U	$O(T_1 T_2 (I_1 + I_2)\log(I_1 + I_2))$	$O(T_1 T_2)$	[52]
less-constrained	O	$O(T_1 T_2 I_1^3I_2^3(I_1 + I_2))$	$O(T_1 T_2 I_1^3I_2^3(I_1 + I_2))$	[29]
less-constrained	U	MAX SNP-hard		[29]
unit-cost	O	$O(w^2\min(T_1 , T_2)\min(L_1, L_2))$	$O(T_1 T_2)$	[41]
1-degree	O	$O(T_1 T_2)$	$O(T_1 T_2)$	[38]

Tree alignment distance

general	O	$O(T_1 T_2 (I_1 + I_2)^2)$	$O(T_1 T_2 (I_1 + I_2)^2)$	[18]
general	U	MAX SNP-hard		[18]
similar	O	$O((T_1 + T_2)\log(T_1 + T_2)(I_1 + I_2)^4s^2)$	$O((T_1 + T_2)\log(T_1 + T_2)(I_1 + I_2)^4s^2)$	[17]

Tree inclusion

general	O	$O(T_1 T_2)$	$O(T_1 \min(D_2L_2))$	[21]
general	O	$O(\Sigma_{T_1} T_2 + m_{T_1}x_2D_2)$	$O(\Sigma_{T_1} T_2 + m_{T_1}x_2)$	[36]
general	O	$O(L_1 T_2)$	$O(L_1\min(D_2L_2))$	[7]
general	U	NP-hard		[22, 32]



Existing Methods

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

Sequencing

Transformation

Experiments

Conclusion

Existing Methods

- PQ-Gram based method
- Binary tree based method
- Filter and Refine



Binary tree based similarity search ^[SIGMOD05]

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

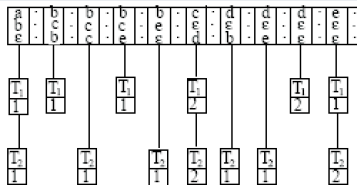
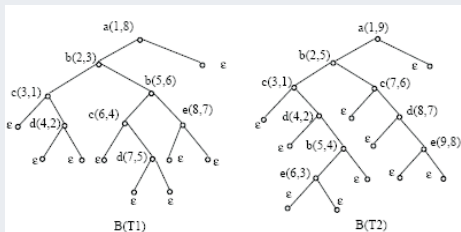
Sequencing

Transformation

Experiments

Conclusion

Binary tree



(a) Inverted File

BRV(T₁) 1 ... 1 ... 0 ... 1 ... 0 ... 2 ... 0 ... 0 ... 2 ... 1 ...

BRV(T₂) 1 ... 0 ... 1 ... 0 ... 1 ... 2 ... 1 ... 1 ... 0 ... 2 ...



pq-gram based similarity search [VLDB05]

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

Sequencing

Transformation

Experiments

Conclusion

pq-gram

Extended Tree T^{pq} :

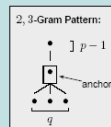
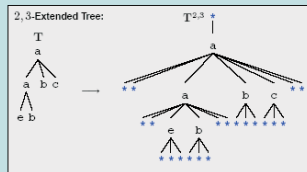
Patch boundaries by adding null nodes (*):

- $\Rightarrow p - 1$ ancestors to the root
- $\Rightarrow q - 1$ nodes before the first and after the last child of each non-leaf node
- $\Rightarrow q$ children to each leaf

pq -Gram G : Subtree of T^{pq} .

- \Rightarrow Anchor node
- \Rightarrow with $p - 1$ ancestors
- \Rightarrow and q children.

Contiguous siblings in G are contiguous siblings in T^{pq} .





pq-gram based Similarity Join [VLDB05]

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

Sequencing

Transformation

Experiments

Conclusion

pq-gram

Extended Tree T^{pq} :

Patch boundaries by adding null nodes (*):

$\Rightarrow p - 1$ ancestors to the root

$\Rightarrow q - 1$ nodes before the first and after the last child of each non-leaf node

$\Rightarrow q$ children to each leaf

pq -Gram G : Subtree of T^{pq} .

\Rightarrow Anchor node

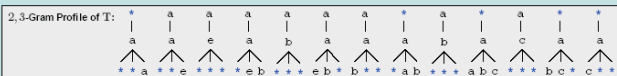
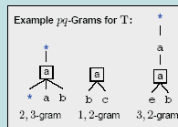
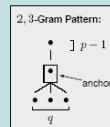
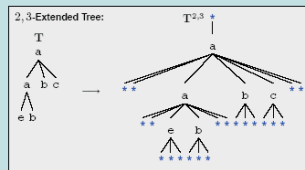
\Rightarrow with $p - 1$ ancestors

\Rightarrow and q children.

Contiguous siblings in G are contiguous siblings in T^{pq} .

pq -gram Profile $P^{p,q}(T)$:

\Rightarrow Bag of all pq -grams of T .





RoadMap

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

Sequencing

Transformation

Experiments

Conclusion



Motivation

- Problem
- Applications
- Similarity Metrics
- Existing Methods



Sequence-Based Tree-Structured Data Similarity Search

- Intuition
- Sequencing
- Edit Distance Transformation



Experiments



Conclusion



Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

Sequencing

Transformation

Experiments

Conclusion

Intuition

- Efficiency
- Both Structural and Textual features

Intuition

- Trees \rightarrow Strings
- Using approximate string search for approximate tree search



Sequence-Based TSearch

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

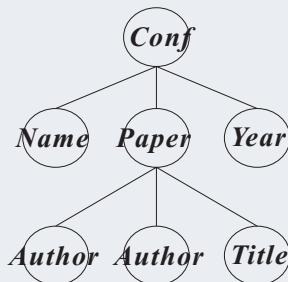
Sequencing

Transformation

Experiments

Conclusion

Sequencing



Name Conf Author Paper Author Paper Title Paper Conf Year Conf



Sequence-Based TSearch

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

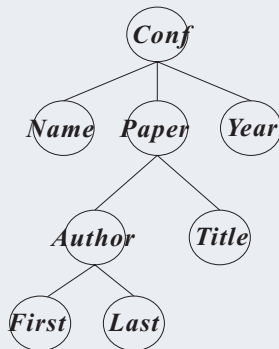
Sequencing

Transformation

Experiments

Conclusion

Sequencing



Name Conf First Author Last Author Paper Title Paper Conf Year Conf



Sequence-Based TSearch

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

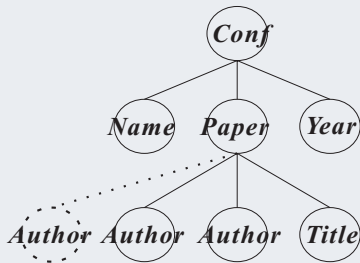
Sequencing

Transformation

Experiments

Conclusion

Insert Leaf Node



Name Conf Author Paper Author Paper Author Paper Title Paper Conf Year Conf



Sequence-Based TSearch

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

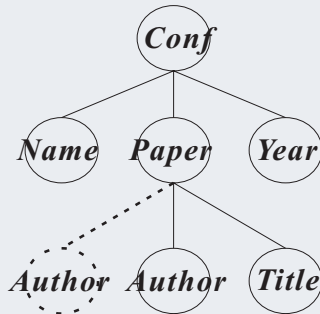
Sequencing

Transformation

Experiments

Conclusion

Delete Leaf Node



Name Conf Author Paper Author Paper Title Paper Conf Year Conf



Sequence-Based TSearch

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

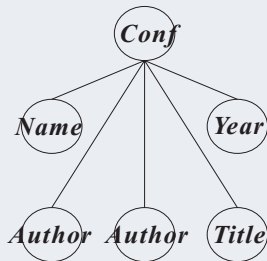
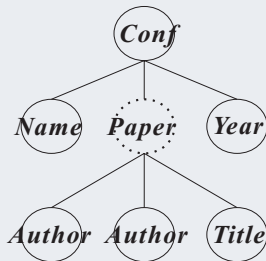
Sequencing

Transformation

Experiments

Conclusion

Delete/Insert Internal Node



Name Conf Author Paper Conf Author Paper Conf Title Paper Conf Year Conf



Edit Distance Transformation

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

Sequencing

Transformation

Experiments

Conclusion

Edit Distance Transformation

$$ted(T, T') \leq ed(S, S') = ed(T, T') \leq C_{max} * ted(T, T')$$



LCST Distance Transformation

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

Sequencing

Transformation

Experiments

Conclusion

LCST Distance Transformation

$$lcstd(T, T') \geq ted(T, T')$$

$$lcstd(T, T') \geq ted(T, T') \geq \frac{1}{c_{max}} * ed(S, S')$$



SCST Distance Transformation

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

Sequencing

Transformation

Experiments

Conclusion

SCST Distance Transformation

$$scstd(T, T') \geq ted(T, T')$$

$$scstd(T, T') \geq ted(T, T') \geq \frac{1}{c_{max}} * ed(S, S')$$



RoadMap

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

Sequencing

Transformation

Experiments

Conclusion



Motivation

- Problem
- Applications
- Similarity Metrics
- Existing Methods



Sequence-Based Tree-Structured Data Similarity Search

- Intuition
- Sequencing
- Edit Distance Transformation



Experiments



Conclusion



Data Sets

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

Sequencing

Transformation

Experiments

Conclusion

Datasets	Average # of elements	Maximal depth	Maximal fan-out
DBLP	13.6	4	8
SIGMOD Record	12.2	4	8
XMark	10.4	6	6
Treebank	11.8	8	6



Pruning Power

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

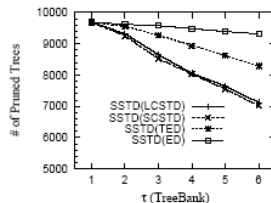
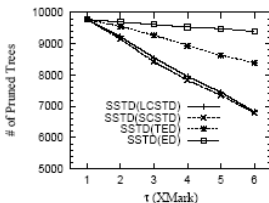
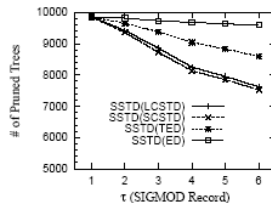
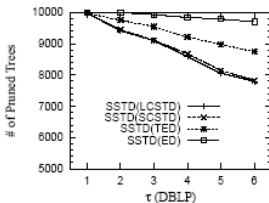
Sequencing

Transformation

Experiments

Conclusion

Average # of Pruned Trees vs. Different Values of τ





Performance

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

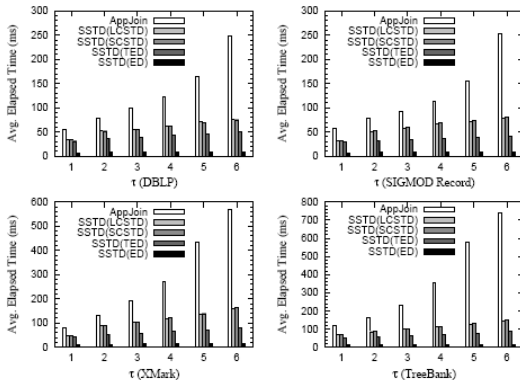
Sequencing

Transformation

Experiments

Conclusion

Average Elapsed Time vs. Different Values of τ





Performance

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

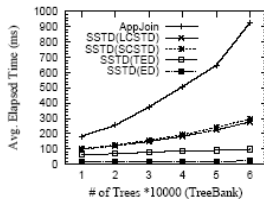
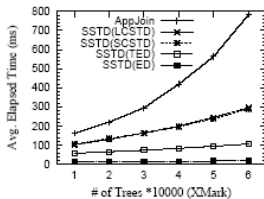
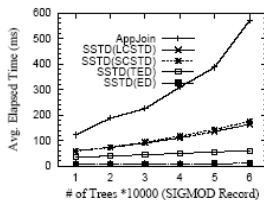
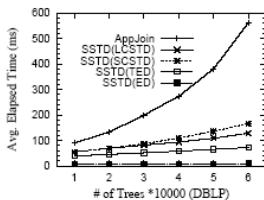
Sequencing

Transformation

Experiments

Conclusion

Average Elapsed Time vs. Different Numbers of Trees ($\tau=3$)





RoadMap

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

Sequencing

Transformation

Experiments

Conclusion



Motivation

- Problem
- Applications
- Similarity Metrics
- Existing Methods



Sequence-Based Tree-Structured Data Similarity Search

- Intuition
- Sequencing
- Edit Distance Transformation



Experiments



Conclusion



Sequence Based Similarity Search

Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

Sequencing

Transformation

Experiments

Conclusion

Sequence Based Similarity Search for Tree-Structured Data

- Trees \rightarrow Strings: Sequencing method
- Using approximate string search for approximate tree search : Edit Distance Transformation
- Efficient methods for different metrics



Tree Search

Guoliang Li

Motivation

Problem

Applications

Metrics

Related Work

TSearch

Intuition

Sequencing

Transformation

Experiments

Conclusion

Thanks!!
Questions?