

# Online Distributed Trace Synchronization

**Masoume Jabbarifar**

**Michel Dagenais**



*December 6, 2012*  
*École Polytechnique, Montreal*

# Outline

**Streaming Mode Incremental Clock Synchronization**

**Reference Node Selection in Dynamic Tree**

**Minimum Spanning Tree Maintenance in Dynamic Tree**

- Key problem
- Methodology
- Experimental results
- Conclusion

# Key problem

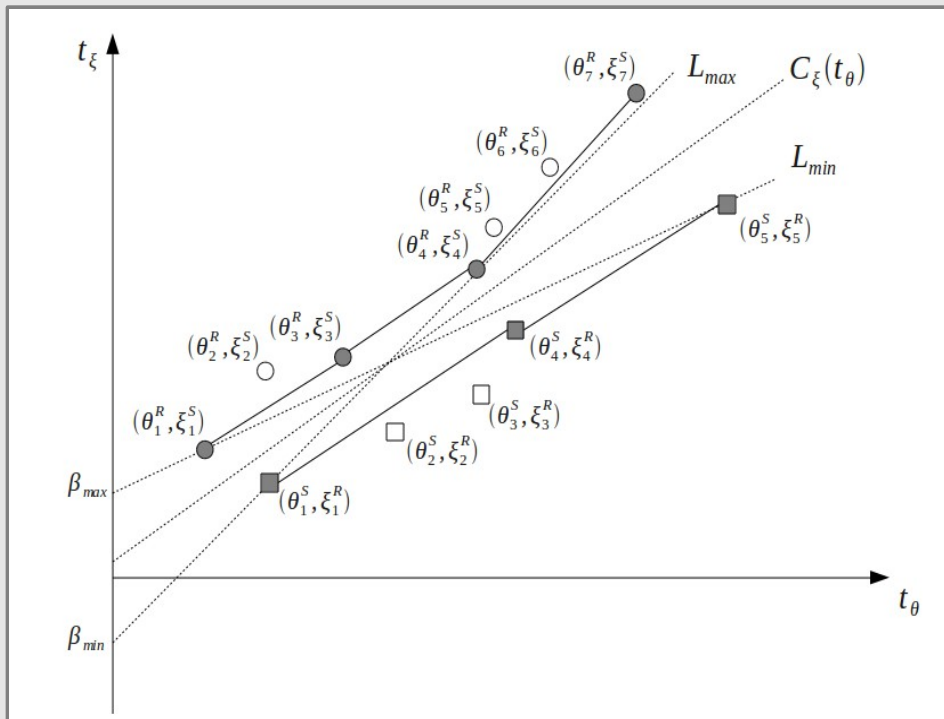
## Problems and solutions

- **Accuracy:** The generated events are at the nanosecond scale.
- **Performance:** The analysis speed must be the same as, or higher than, the data rate.
- **Data buffering:** The gathered data from a computer cluster is huge but the required data is smaller.
- **Delay:** Little delay must be added by analysis for tracing to diagnose problems and attacks in real time
- **Convex-hull algorithm:** This algorithm guarantees the best accuracy.
- **Incremental algorithm:** It refreshes synchronization as it receives accurate data.
- **Layered improvements:** Proposed online synchronization methods improve performance in different layers;
  - 1) Individual connection
  - 2) In a computer network
  - 3) Time reference updates

# **“Streaming Mode Incremental Clock Synchronization”**

# Two Clocks Synchronization

## Convex-Hull



- Sent and Received messages sets
- Guarantees no message inversion
- Two lines with Max & Min slope

$$L_{max} = \alpha_{max} \theta + \beta_{min}$$

$$L_{min} = \alpha_{min} \theta + \beta_{max}$$

$$Accuracy = \alpha_{max} - \alpha_{min}$$

- The bisector of the angle formed by these two lines

$$C_A(t) = \alpha C_B(t) + \beta$$

# Time Interval based approaches (1)

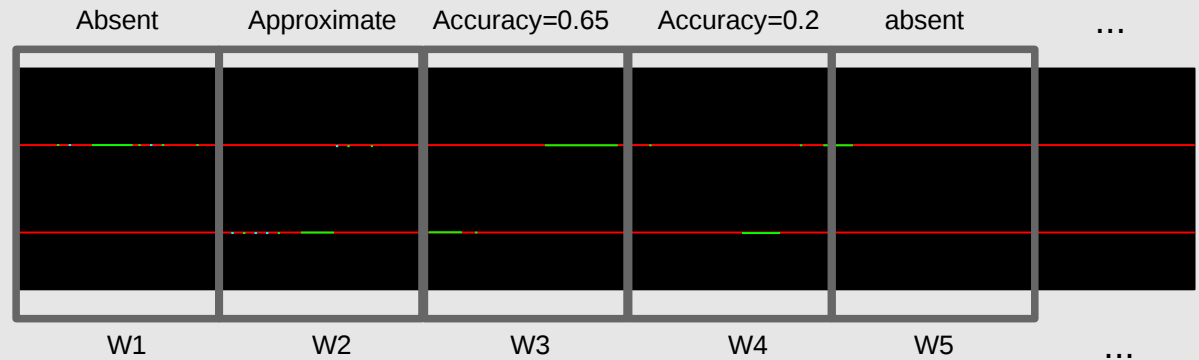
## Independent Window

### Advantages

- Performance
- No buffering
- Simple to implement

### Disadvantages

- Accuracy

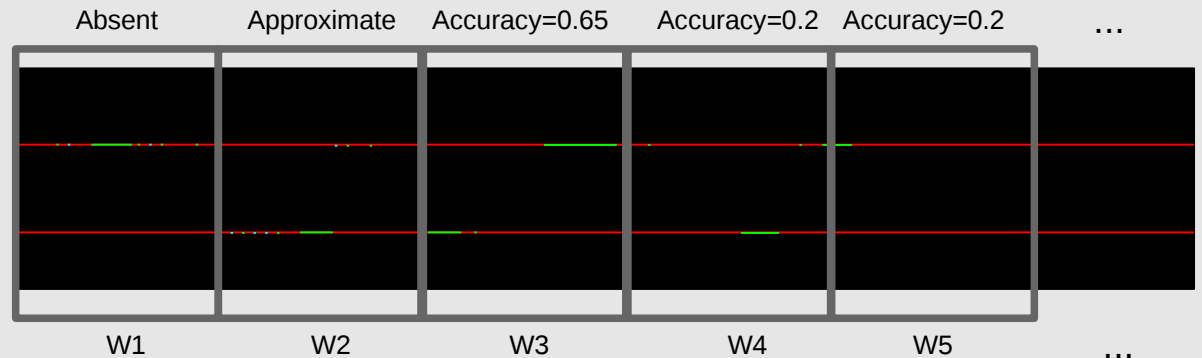


# Time Interval based approaches (2)

## Replacement approach

### Advantages

- Performance
- Simple to implement

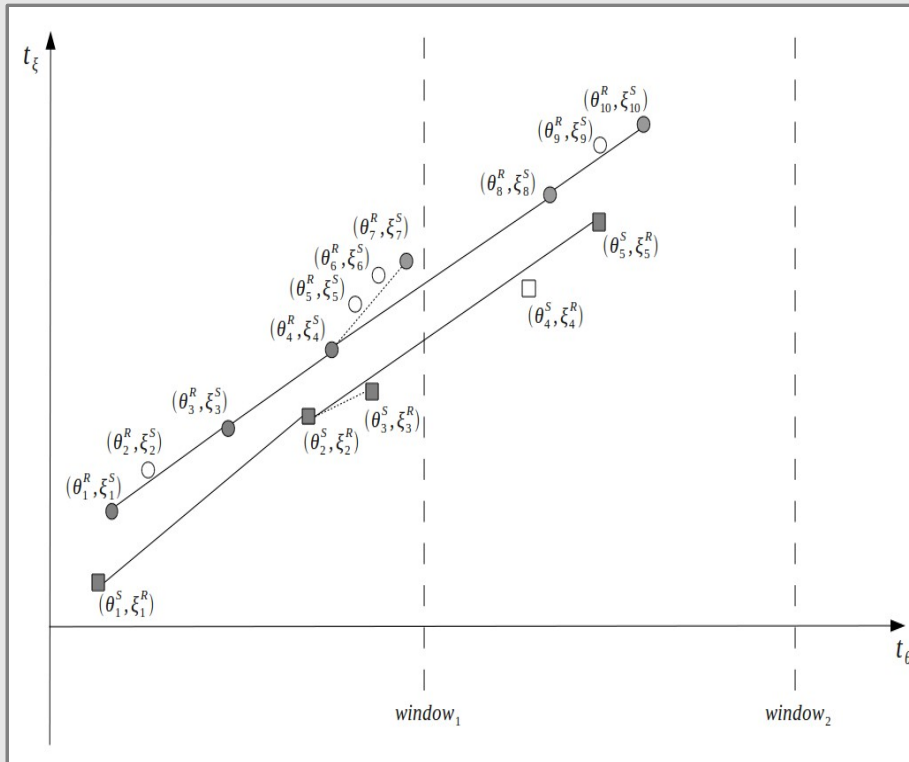


### Disadvantages

- High level accuracy is unobtainable

# Time Interval based approaches (3)

## Correlated approach



### Advantages

- The highest level of accuracy
- No buffering

### Disadvantages

- Processing is postponed to the end of each window



# Fully incremental Approach (1)

## No window concept

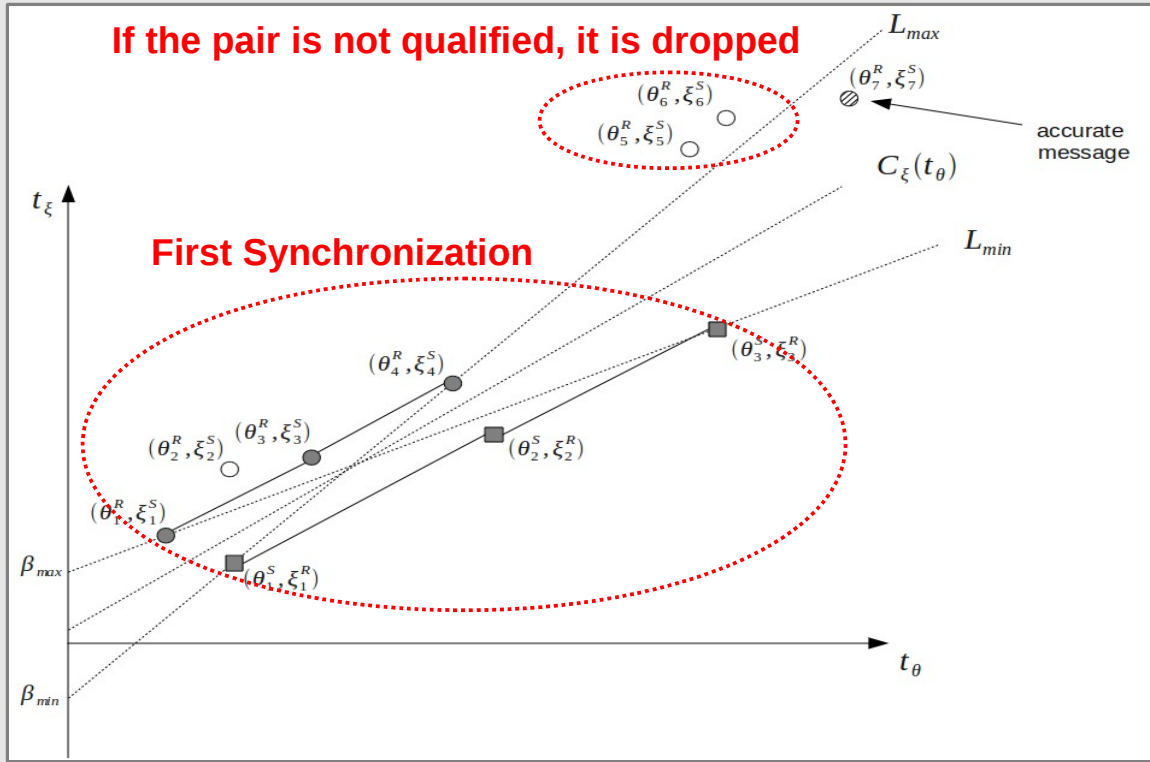
- The highest level of accuracy
- No buffering
- No delay



# Fully Incremental Approach (2)

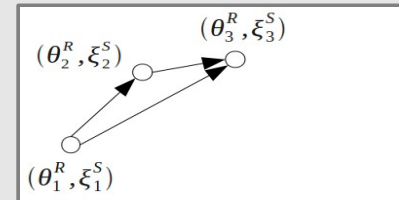
Step1: wait to establish at least three messages in each direction

Step2: Look for accurate messages

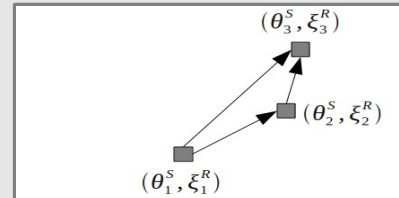


Detection method:

Lower bound



Upper bound

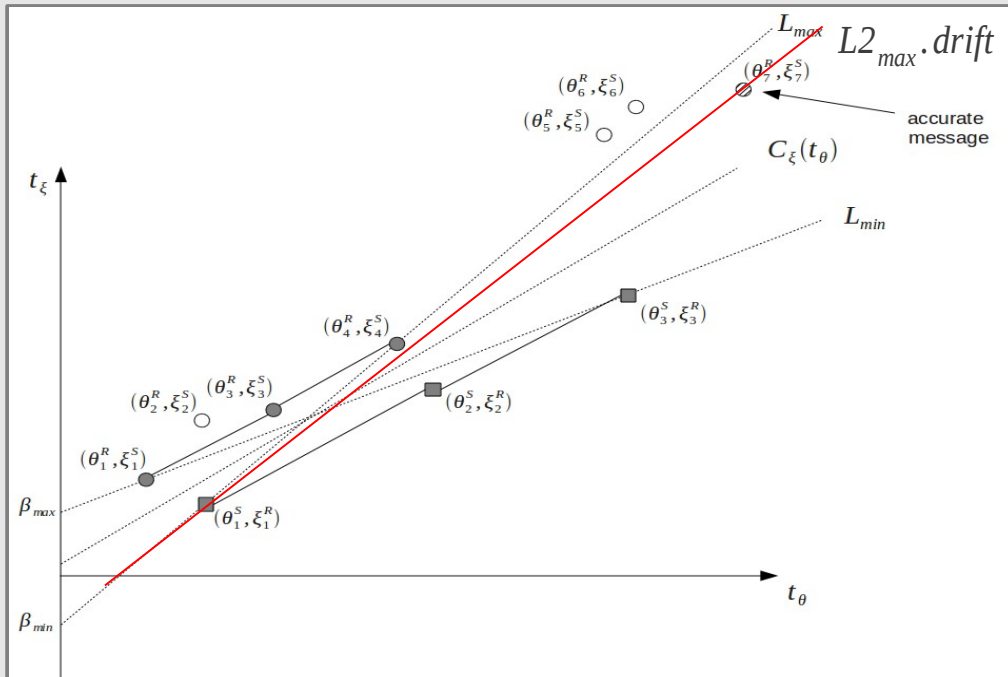


# Fully Incremental Approach (3)

## Accurate Message Improves Synchronization Accuracy

- accuracy is the difference between the minimum and maximum possible drifts between the two clocks

$$Accuracy1 = L1_{max} \cdot drift - L1_{min} \cdot drift$$



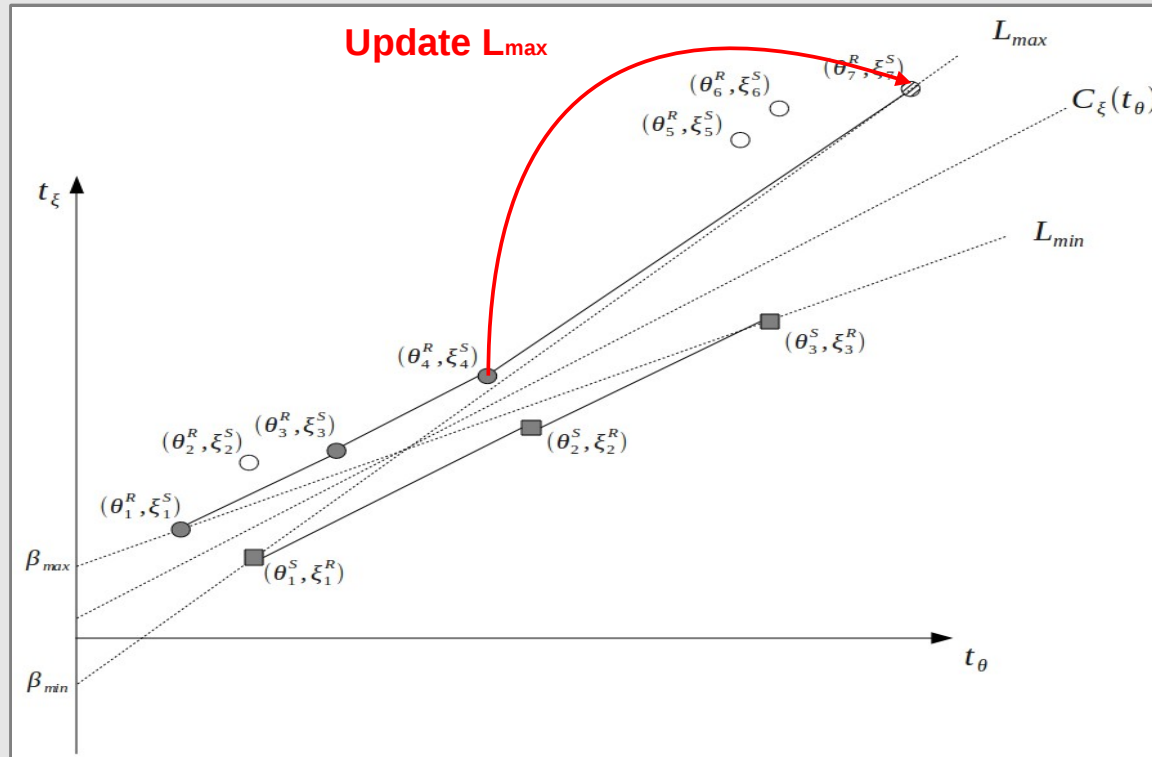
$$L1_{max} \cdot drift > L2_{max} \cdot drift$$

$$Accuracy1 > Accuracy2$$

# Fully Incremental Approach (4)

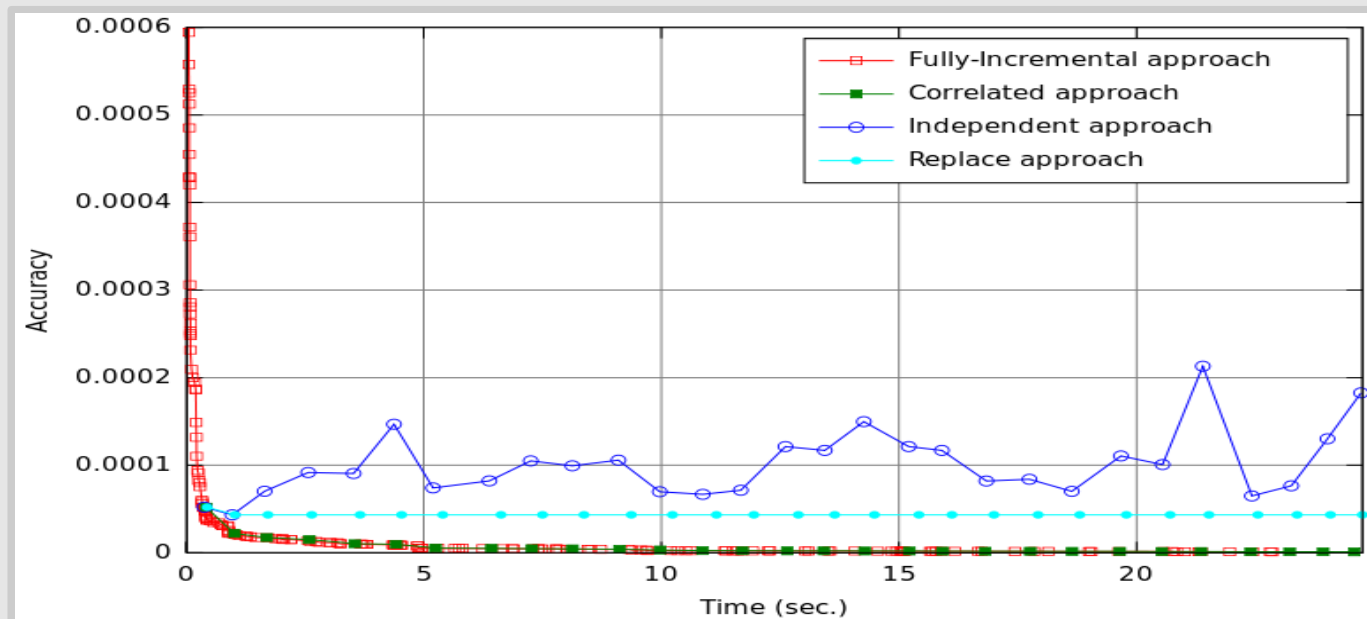
## Step3: adjust-bounds

- Only  $L_{max}$  or  $L_{min}$  is changed when an accurate pair is received

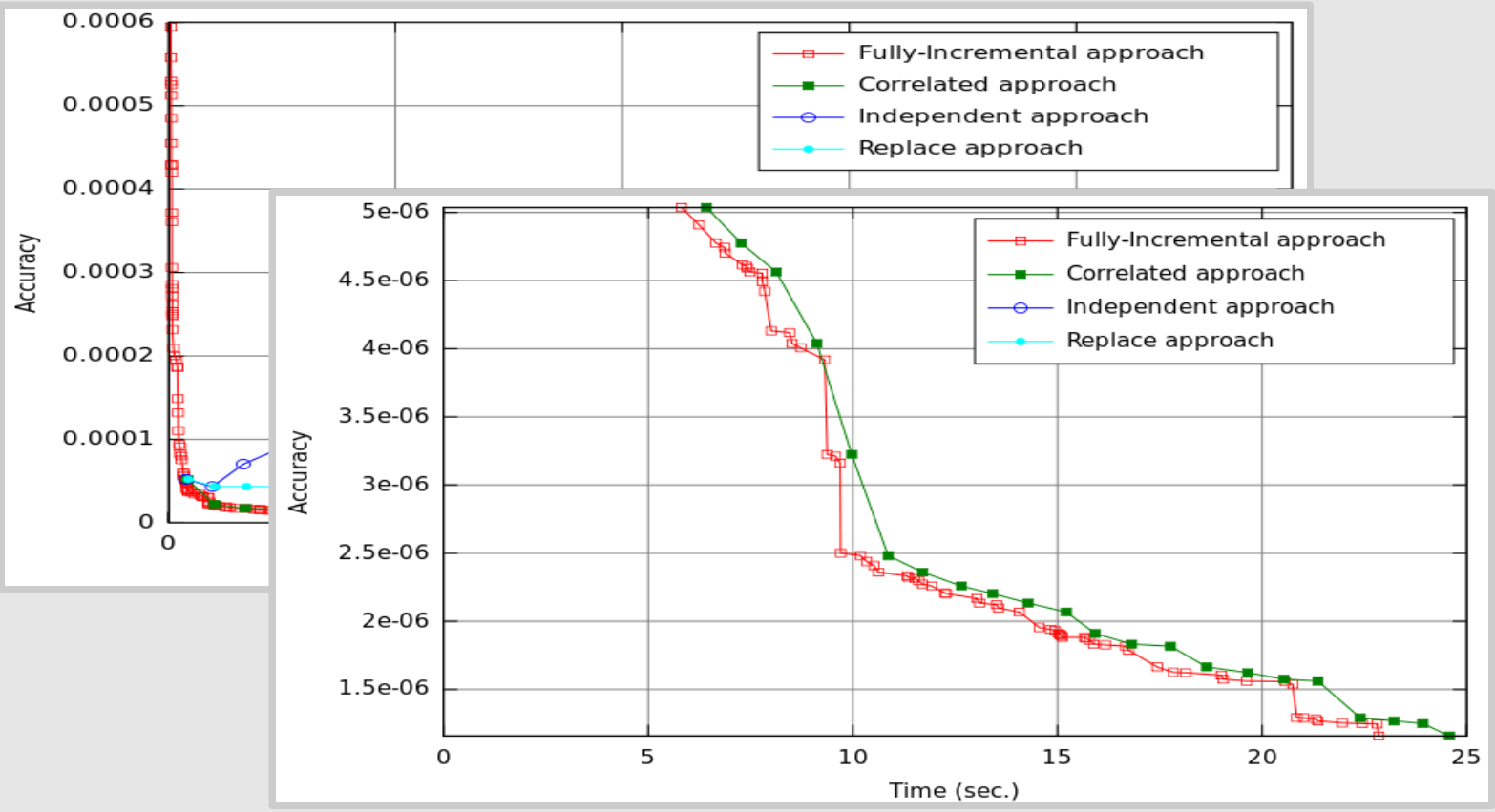


Each synchronization in the Fully Incremental approach requires  $O(1)$  time, on average

# Results (1)



# Results (2)



Average delay = 0.91 sec

# Fully incremental Approach (5)

## No window consideration

### Advantages

- Performance: Each synchronization requires  $O(1)$  time, on average.
- High level time accuracy
- No delay
- No buffering

### Features

- **Analysis:** Appropriate data are filtered prior to synchronization computations.

# **“Reference Node Selection in Dynamic Tree“**



# Key problem

## Reference Node Selection

- 1) For any node in Spanning Tree, the shortest paths to every other node are computed.
  - 2) The best time reference node (RN) is the node which has the smallest paths to all.
  - 3) All nodes in the network synchronize their time with the reference node through those paths.
- It takes  $O(n^2)$  time to find the RN.
  - A fixed RN?
    - Costs synchronization accuracy.
    - A single point of failure.

# Reference Node Position

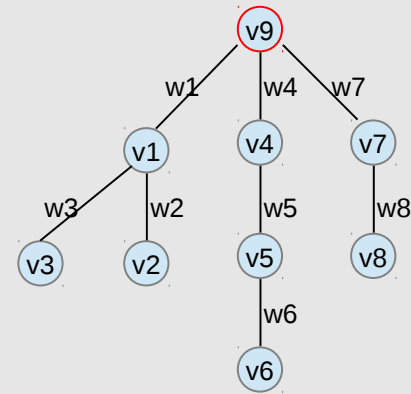
## Total accuracy

- The position of RN is critical to decrease the total time conversion error through all paths.

$$T: \langle v_1, v_2, \dots, v_n \rangle$$

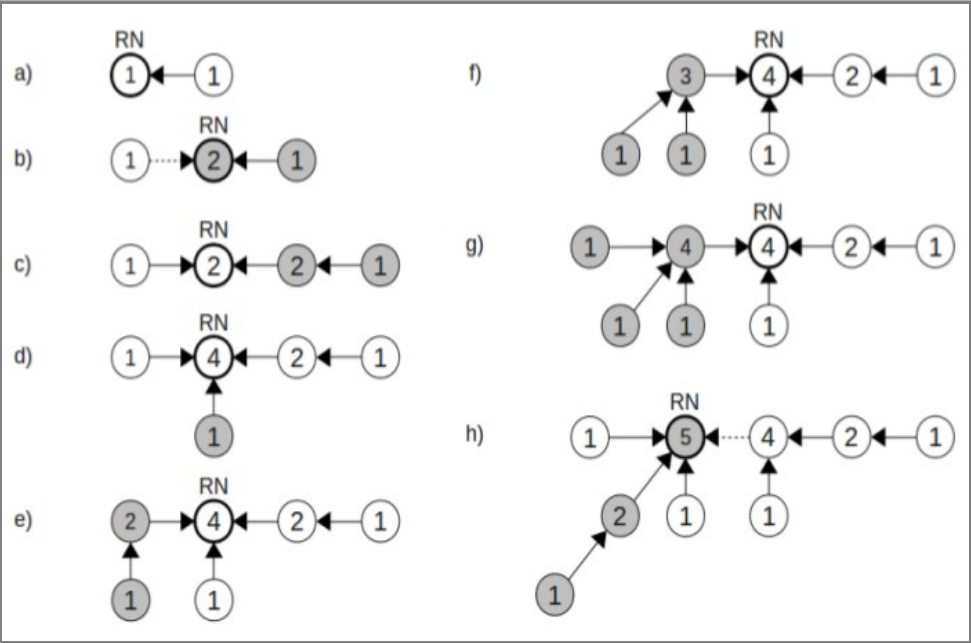
$$P_{v_i} \langle v_i, v_{i+1}, \dots, RN \rangle$$

$$\text{Total Synchronization Error}_T = \sum_{i=1}^n \sum_{j=1}^l \text{weight}_{e_j} \text{ on } P_{v_i}$$



# Dynamic Reference Node (1)

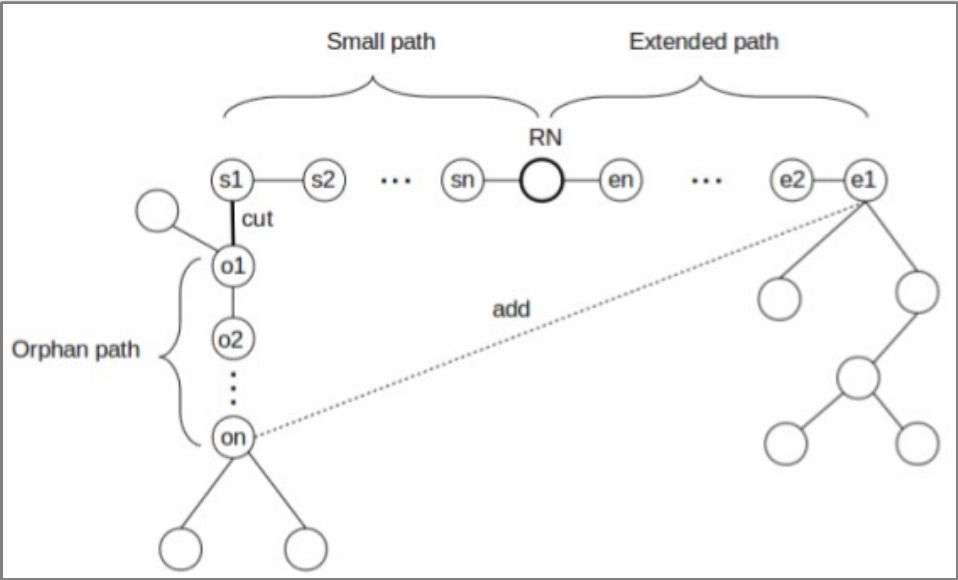
## New Node Connection



- The number inside each vertex shows the Descendant Size
- Propagation along the path from the parent of “v” to the RN
- Comparison between two nodes

# Dynamic Reference Node (2)

## Cycle/ Cut and Add (1)

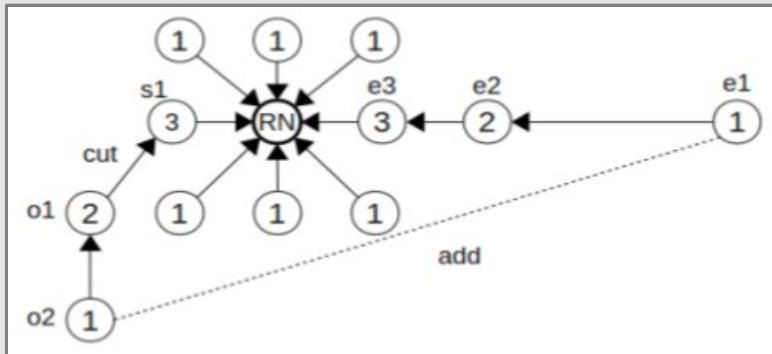
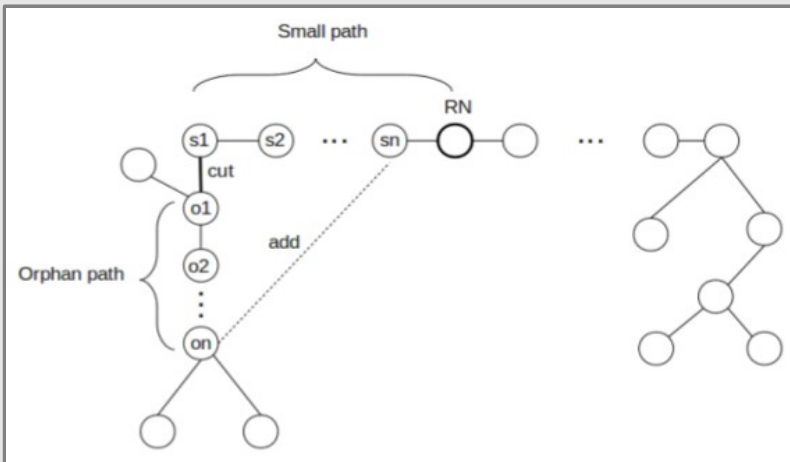


RN is in Extended path

- Orphan path:  $\langle o_2, \dots, o_{n-1} \rangle$
- Small path:  $\langle s_1, s_2, \dots, s_n \rangle$
- Extended path:  $\langle e_1, e_2, \dots, e_n \rangle$

# Dynamic Reference Node (3)

## Cycle/ Cut and Add (2)



- RN does not change when:
  - Cut and add are in one side of RN
  - RN has many branches

$reverse(e_1)$

$$\Delta = \underbrace{\xi(o_1)}_{\text{cut impact}}$$

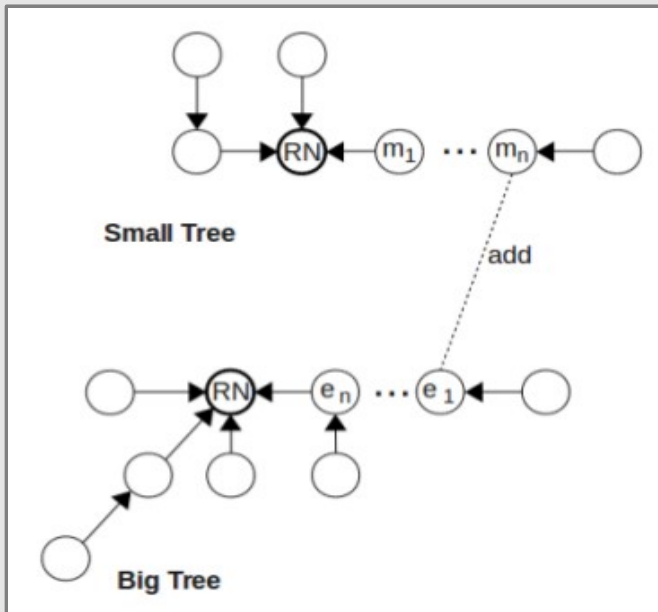
the new DescendantSize for previous reference node

$$\xi(\chi) = \overbrace{treeSize(\chi) - [\xi(parent(\chi)) + \Delta]}^{e_n}$$

$$RN(\tau_i) = \begin{cases} \chi & \text{if } \xi(\chi) > \underbrace{treeSize(\chi)/2}_{\text{tree balance value}} \\ \text{search in extendedPath} & \text{otherwise} \end{cases}$$

# Dynamic Reference Node (4)

## Joining Two Trees



- 1) Edge  $(RN_{ST}, RN_{BT}) \rightarrow RN_{BT}$
- 2) Edge  $(RN_{ST}, e_i) \rightarrow$   
Reverse path  $\langle e_i, e_{i+1}, \dots, RN_{BT} \rangle$
- 3) Edge  $(m_i, RN_{BT}) \rightarrow RN_{BT}$   
Reverse path  $\langle RN_{ST}, m_1, \dots, m_n \rangle$
- 4) Edge  $(m_{ST}, e_{BT}) \rightarrow RN_{BT}$   
Reverse path  $\langle RN_{ST}, m_1, \dots, m_n \rangle$   
Candidate  $RN \in \langle e_i, e_{i+1}, \dots, RN_{BT} \rangle$

# Data set (1)

## Number of each operation, from one million operations

	Nodes	Insertion	Join	Cycle		updateEdge
				Stay <sup>1</sup>	Remove <sup>2</sup>	
<i>Dataset</i> <sub>1</sub>	10000	4991	2503	45449	946879	178
<i>Dataset</i> <sub>2</sub>	20000	9892	5052	76404	908556	96
<i>Dataset</i> <sub>3</sub>	30000	15005	7496	102672	874761	66
<i>Dataset</i> <sub>4</sub>	40000	19955	10021	125650	844322	52
<i>Dataset</i> <sub>5</sub>	50000	24959	12519	145733	816753	36
<i>Dataset</i> <sub>6</sub>	60000	29953	15022	164104	790885	36

<sup>1</sup> the new connection stay in loop and other edge in cycle is removed by MST algorithm.

<sup>2</sup> the new connection has the highest weight in cycle and is removed by MST algorithm

## RN changes

	Insertion					Join					Cycle				
	Number <sup>1</sup>	<i>win</i> <sub>RN</sub> <sup>2</sup>	%	<i>lose</i> <sub>RN</sub> <sup>3</sup>	%	Number <sup>4</sup>	<i>win</i> <sub>RN</sub>	%	<i>lose</i> <sub>RN</sub>	%	Number <sup>5</sup>	<i>win</i> <sub>RN</sub>	%	<i>lose</i> <sub>RN</sub>	%
<i>Dataset</i> <sub>1</sub>	4991	732	15%	4259	85%	2503	1464	58%	1039	42%	45449	1186	3%	44263	97%
<i>Dataset</i> <sub>2</sub>	9892	1435	14%	8457	86%	5052	2972	70%	2080	30%	76404	1259	2%	75145	98%
<i>Dataset</i> <sub>3</sub>	15005	2270	15%	12735	85%	7496	4439	59%	3057	41%	102672	1508	2%	101164	98%
<i>Dataset</i> <sub>4</sub>	19955	2986	15%	16969	85%	10021	5847	58%	4174	42%	125650	1365	1%	124285	99%
<i>Dataset</i> <sub>5</sub>	24959	3675	15%	21284	85%	12519	7298	58%	5221	42%	145733	1432	1%	144301	99%
<i>Dataset</i> <sub>6</sub>	29953	4373	15%	25580	85%	15022	8839	59%	6183	41%	164104	1776	1%	162328	99%

<sup>1</sup> The total number of cases where a vertex add to existent tree. Note that other cases belong to two new vertices connections which makes a new tree. In recent case one of vertices is selected as RN and there is no computation to find RN.

<sup>2</sup> The number of cases that RN changes.

<sup>3</sup> The number of cases that RN does not change.

<sup>4</sup> The number of cases that two trees merge in the forest.

<sup>5</sup> The number of cases that an edge makes cycle in one of trees in the forest.

## Data sets (2)

### Merge positions

	$RN_s - RN_b$	$m_i - RN_b$	$RN_s - e_i$	$m_i - e_i$
<i>Dataset</i> <sub>1</sub>	0	190	80	2233
<i>Dataset</i> <sub>2</sub>	0	440	161	4451
<i>Dataset</i> <sub>3</sub>	0	595	232	6669
<i>Dataset</i> <sub>4</sub>	0	819	339	8863
<i>Dataset</i> <sub>5</sub>	0	1035	420	11064
<i>Dataset</i> <sub>6</sub>	0	1238	474	13310

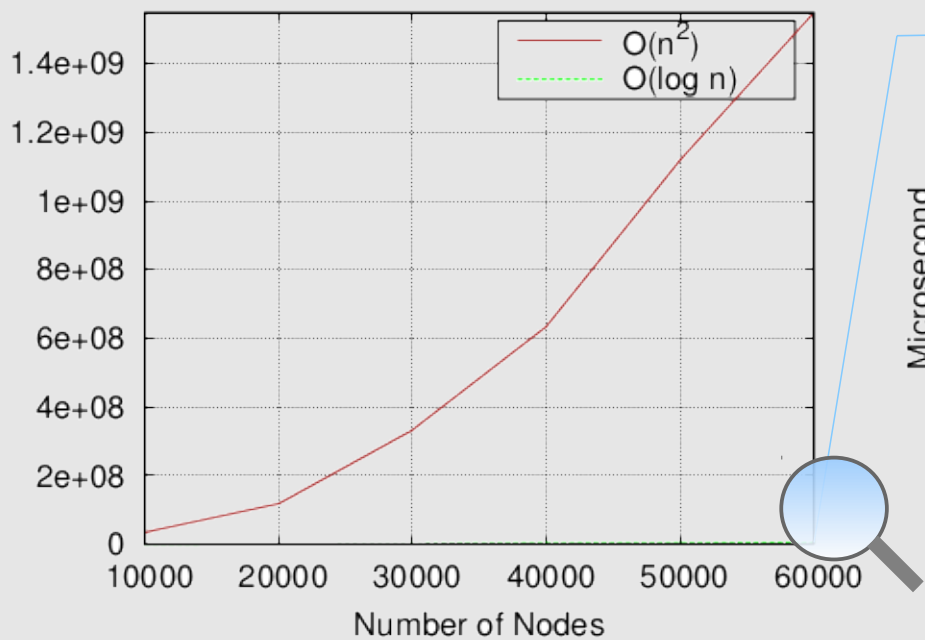
### Number of descendant Size update in each operation

	Insertion			Join			Cycle		
	Number	Update	Avg. Distance to RN	Number	Update	Avg. Distance to RN	Number	Update	Avg. Distance to RN
<i>Dataset</i> <sub>1</sub>	4991	89985	18	2503	50930	20	45449	3431773	75.5
<i>Dataset</i> <sub>2</sub>	9892	246827	25	5052	128818	25.5	76404	7614996	100
<i>Dataset</i> <sub>3</sub>	15005	366967	24	7496	191811	26	102672	12250493	119
<i>Dataset</i> <sub>4</sub>	19955	513153	26	10021	275419	27.5	125650	15740773	127
<i>Dataset</i> <sub>5</sub>	24959	762460	31	12519	398419	32	145733	18804027	129
<i>Dataset</i> <sub>6</sub>	29953	1018792	34	15022	540992	36	164104	23822448	145



# Result

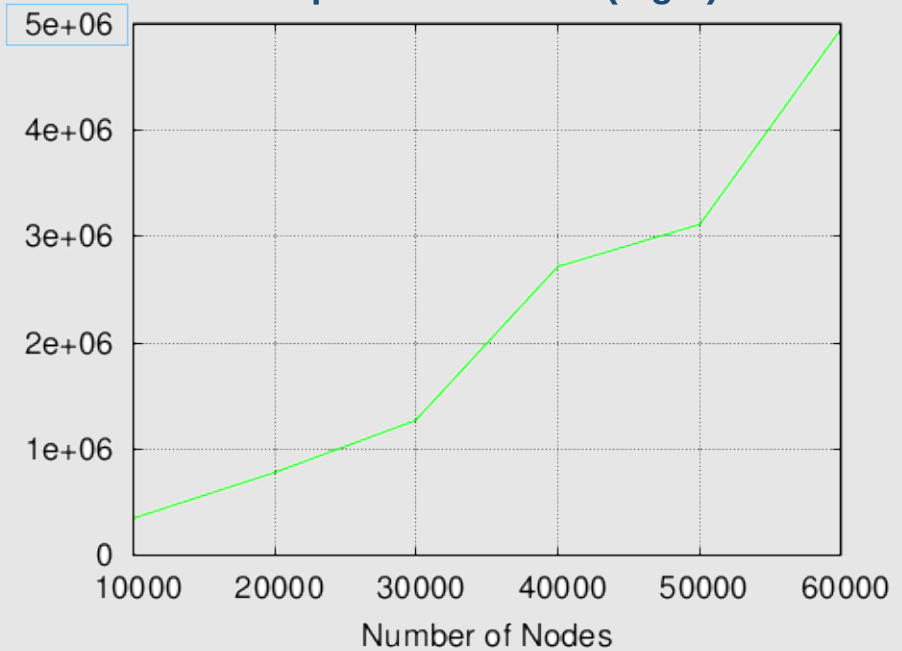
Updating the reference node in a dynamic network with one million operations



10000 nodes= 32.53 sec.

60000 nodes= 1549.92 sec.

## Proposed method $O(\log n)$



10000 nodes= 0.33 sec.

60000 nodes= 4.92 sec.

# **“Minimum Spanning Tree Maintenance in Dynamic Tree“**

# Kruskal's Algorithm for Minimum Spanning Tree

- Running Time =  $O(m \log n)$   $m$  = edges,  $n$  = nodes
- The steps are:
  - 1) The edges are placed in a priority queue
  - 2) Until we have added  $n-1$  edges
    - i. Extract the lowest edge from the queue
    - ii. If it forms a cycle, eliminate it
    - iii. Else add it to the tree

# Challenges in Dynamic Minimum Spanning Tree

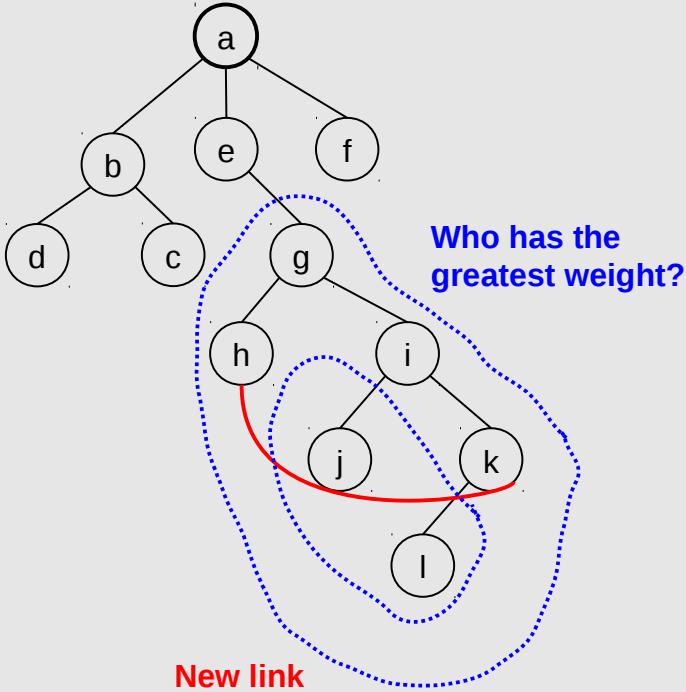
- Two cases can occur:
  - 1) Edge is in Tree and new weight is less than previous (no effect)
    - It is the characteristic of Fully incremental Approach

$$Accuracy_{new} < Accuracy_{previous}$$

- 1) Edge is not in Tree
  - Add without cycle (no effect)
  - Add with cycle (effect)

# Dynamic MST

Loop in MST

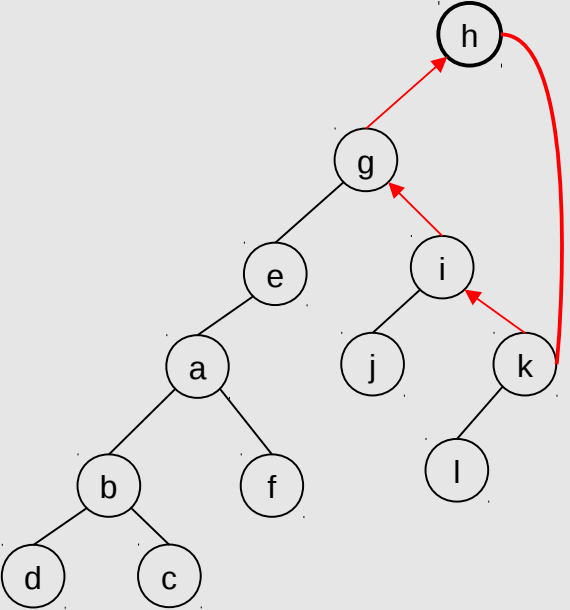


Evert (h)

→

parent(g)= h  
parent(e)= g  
parent(a)= e

Make vertex h the root of its tree



# Data set

## Number of each operation, from one million operations

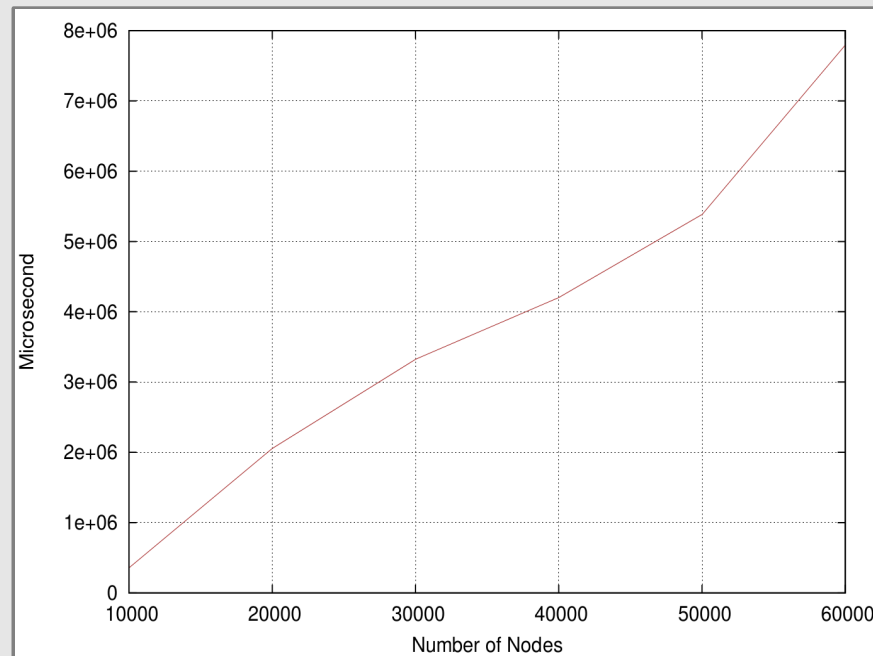
	Nodes	Added link <sub>1</sub>	Removed link <sub>2</sub>	Total
<i>Dataset1</i>	10000	45449	946879	992328
<i>Dataset2</i>	20000	76404	908556	984960
<i>Dataset3</i>	30000	102672	874761	977433
<i>Dataset4</i>	40000	125650	844322	969972
<i>Dataset5</i>	50000	145733	816753	962486
<i>Dataset6</i>	60000	164104	790885	954989

- 1) The new connection is added to the MST and one of other edges in the cycle is removed by MST algorithm
- 2) The new connection has the highest weight in cycle and is removed by MST algorithm

# Result

- Old approach requires 1.440216 sec. to find MST once in a cluster with 10000 simulated nodes.
- For 992328 changes in MST:  $992328 * 1.440216 = 1429166.662848$  sec.  $\sim$  396 hours = 16.5 days!!

## Proposed method $O(\log n)$



10000 nodes= 0.36 sec.

60000 nodes= 7.79 sec.

# Conclusion

- Resynchronization rate per edge
  - Performance: Each synchronization requires  $O(1)$  time, on average.
  - High level time accuracy
  - No delay
  - No buffering
- Resynchronization rate per network
  - Reference node selection
  - Synchronization path update
  - Performance: Each update requires  $O(\log n)$  time, on average.
- Error reduction and continuity
- Scalability
- Robustness



Thank you

# DORSAL

---

Online Distributed Trace Synchronization

[www.lttng.org](http://www.lttng.org)

E-Mail:

[masoume.jabbarifar@polymtl.ca](mailto:masoume.jabbarifar@polymtl.ca)

[michel.dagenais@polymtl.ca](mailto:michel.dagenais@polymtl.ca)