# A Tree Clock Data Structure for Causal Orderings in Concurrent Executions

Umang Mathur
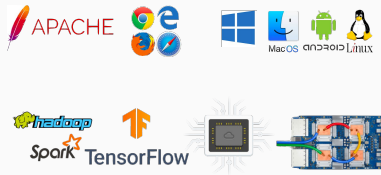
Hünkar Can Tunç

Andreas Pavlogiannis

Mahesh Viswanathan
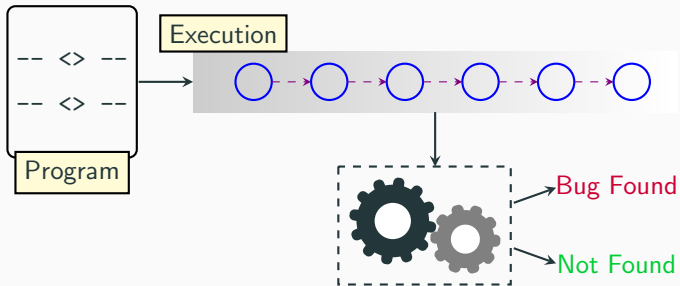
NUS
National University
of Singapore

AARHUS UNIVERSITY

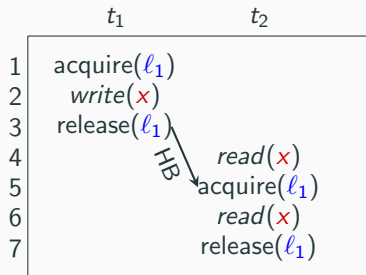UNIVERSITY OF
ILLINOIS
URBANA-CHAMPAIGN

- Ubiquitous computing paradigm.
- Analysis of concurrent programs is a major challenge.
- We need more **efficient** algorithms and data structures.

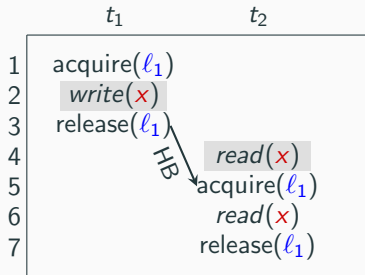# Dynamic Analyses for Detecting Concurrency Bugs



- Widely adopted (e.g., ThreadSanitizer, Helgrind).
- Requires establishing a causal ordering between the events.
- Causality is typically represented as a partial order.

# Happens-Before (HB) Partial Order

# Happens-Before (HB) Partial Order



|   | $t_1$ | $t_2$ |
|---|-------|-------|
| 1 | acquire($\ell_1$) | |
| 2 | write($x$) | |
| 3 | release($\ell_1$) | |
| 4 | | read($x$) |
| 5 | | acquire($\ell_1$) |
| 6 | | read($x$) |
| 7 | | release($\ell_1$) |

HB

$\rightarrow$ Events $e_2$ and $e_4$ are concurrent.

## Happens-Before (HB) Partial Order



|   | $t_1$ | $t_2$ |
|---|-------|-------|
| 1 | acquire($\ell_1$) | |
| 2 | write($x$) | |
| 3 | release($\ell_1$) | |
| 4 | | read($x$) |
| 5 | | acquire($\ell_1$) |
| 6 | | read($x$) |
| 7 | | release($\ell_1$) |

$\rightarrow$ Events $e_2$ and $e_4$ are concurrent.

$\rightarrow$ Events $e_2$ and $e_6$ are **not** concurrent.

## Happens-Before (HB) Partial Order



$\rightarrow$ Events $e_2$ and $e_4$ are concurrent.

$\rightarrow$ Events $e_2$ and $e_6$ are **not** concurrent.

Happens-Before defines data races in various memory models.

Tree Clocks: A new data structure

- Can be used to compute Happens-Before efficiently.
    - Optimal data structure for Happens-Before.

Tree Clocks: A new data structure

- Can be used to compute Happens-Before efficiently.
    - Optimal data structure for Happens-Before.
- Versatile data structure.
    - Other partial orders can also be computed efficiently.
    - Schedulable-Happens-Before
    - Mazurkiewicz

Tree Clocks: A new data structure

- Can be used to compute Happens-Before efficiently.
    - Optimal data structure for Happens-Before.
- Versatile data structure.
    - Other partial orders can also be computed efficiently.
    - Schedulable-Happens-Before
    - Mazurkiewicz
- Significant speedups compared to vector clocks.

## Background: Vector Timestamps

- The knowledge set of a thread $t$ can be succinctly captured by a function:

$$V_t : \text{Threads} \to \mathbb{N}$$

- $V_t(t')$ gives the last event of $t'$ that $t$ knows about.
- $t$ knows about all preceding events as well.

## Background: Vector Timestamps

- The knowledge set of a thread $t$ can be succinctly captured by a function:

$$V_t : \text{Threads} \rightarrow \mathbb{N}$$

- $V_t(t')$ gives the last event of $t'$ that $t$ knows about.
- $t$ knows about all preceding events as well.

$$V_{t_2} = [\overset{t_1}{27}, \overset{t_2}{3}, \overset{t_3}{9}, \overset{t_4}{45}, \overset{t_5}{17}, \overset{t_6}{26}]$$

- $t_2$ knows of the first $27$ events of $t_1$.
- $t_2$ has performed $3$ events.

## Background: Vector Timestamps

- The knowledge set of a thread $t$ can be succinctly captured by a function:

$$V_t : \text{Threads} \to \mathbb{N}$$

- $V_t(t')$ gives the last event of $t'$ that $t$ knows about.
- $t$ knows about all preceding events as well.

$$V_{t_2} = [\overset{t_1}{27}, \overset{t_2}{3}, \overset{t_3}{9}, \overset{t_4}{45}, \overset{t_5}{17}, \overset{t_6}{26}]$$

- $t_2$ knows of the first $27$ events of $t_1$.
- $t_2$ has performed $3$ events.

### Operations

$$
\begin{aligned}
V_1 \sqsubseteq V_2 \quad &\text{iff} \quad \forall t \colon V_1(t) \leq V_2(t) \qquad &&\text{(Comparison)} \\
V_1 \sqcup V_2 \quad &= \quad \lambda t \colon \max(V_1(t), V_2(t)) \qquad &&\text{(Join)}
\end{aligned}
$$

## Background: Implementing Vector Timestamps

Just use a vector clock $\qquad VC_t = [27, 3, 9, 45, 17, 26]$

**Vector Clock Join** $VC_1 \leftarrow VC_1 \sqcup VC_2$

- For each thread $t$:
  - If $VC_1[t] < VC_2[t]$
    - $VC_1[t] \leftarrow VC_2[t]$

**Vector Clock Copy** $VC_1 \leftarrow VC_2$

- For each thread $t$:
  - $VC_1[t] \leftarrow VC_2[t]$

Each operation takes $O(\mathcal{T})$ time, for $\mathcal{T}$ threads

## Background: Computing Happens-Before with Vector Clocks

- One vector clock $\mathbb{C}_t$ per thread $t$
- One vector clock $\mathbb{C}_\ell$ per lock $\ell$

---

**Algorithm:** Happens-Before (HB)

---

1 **procedure** acquire($t$, $\ell$)

2   |   $\mathbb{C}_t \leftarrow \mathbb{C}_t \sqcup \mathbb{C}_\ell$;   /* Vector clock join */

3 **procedure** release($t$, $\ell$)

4   |   $\mathbb{C}_\ell = \mathbb{C}_t$;         /* Vector clock copy */

---

## Background: Computing Happens-Before with Vector Clocks

- One vector clock $\mathbb{C}_t$ per thread $t$
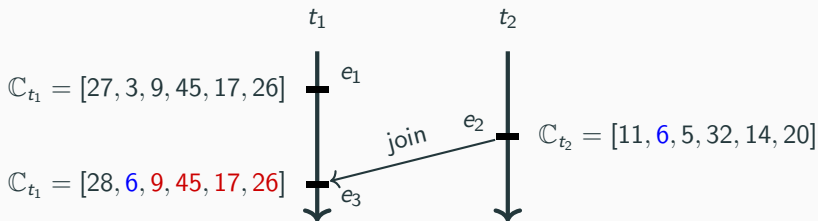- One vector clock $\mathbb{C}_\ell$ per lock $\ell$

---

**Algorithm:** Happens-Before (HB)

---

1 **procedure** acquire$(t, \ell)$

2 $\quad$ $\mathbb{C}_t \leftarrow \mathbb{C}_t \sqcup \mathbb{C}_\ell$;  /* Vector clock join */

3 **procedure** release$(t, \ell)$

4 $\quad$ $\mathbb{C}_\ell = \mathbb{C}_t$;  /* Vector clock copy */

---

- Every vector clock operation costs $O(\mathcal{T})$
  - $\mathcal{T}$ is the number of threads
- When threads are many, the complexity is quadratic $O(\mathcal{N} \cdot \mathcal{T})$
  - $\mathcal{N}$ is the number of acquire/release events

## Overhead of Vector Clocks

- Every vector clock join takes $O(\mathcal{T})$ time.
- Certain steps in the join operation can be vacuous.

$$\mathbb{C}_{t_1} = [27, 3, 9, 45, 17, 26]$$

$$\mathbb{C}_{t_1} = [28, 6, 9, 45, 17, 26]$$

$t_1$      $t_2$

$e_1$

join   $e_2$   $\mathbb{C}_{t_2} = [11, 6, 5, 32, 14, 20]$
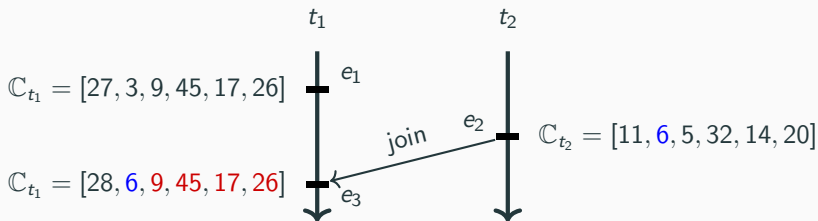
$e_3$

8

## Overhead of Vector Clocks

- Every vector clock join takes $O(\mathcal{T})$ time.
- Certain steps in the join operation can be vacuous.

### Can we do sub-linear joins?

$\rightarrow$ Sub-linear means skip looking at certain entries. How?

$\rightarrow$ Tree clocks address this challenge.

$\mathbb{C}_{t_1} = [27, 3, 9, 45, 17, 26]$

$\mathbb{C}_{t_1} = [28, 6, 9, 45, 17, 26]$

$\mathbb{C}_{t_2} = [11, 6, 5, 32, 14, 20]$

$t_1$      $t_2$

$e_1$

join      $e_2$

$e_3$

## Our Contribution: Tree Clock Data Structure

- **Drop-in** replacement of vector clocks.
- Tree clocks maintain information hierarchically.
    - Nodes store local times of a thread + metadata.
    - Tree structure records how information has been obtained transitively.

- **Drop-in** replacement of vector clocks.
- Tree clocks maintain information hierarchically.
  - Nodes store local times of a thread + metadata.
  - Tree structure records how information has been obtained transitively.

$$\mathbb{C}_{t_4} = [\overset{t_1}{1}, \overset{t_2}{2}, \overset{t_3}{13}, \overset{t_4}{9}] \qquad \longrightarrow$$

$t_4, 9, \perp$

$t_3, 13, 5 \quad t_2, 2, 1$

$t_1, 1, 1$

## Our Contribution: Tree Clock Data Structure

- **Drop-in** replacement of vector clocks.
- Tree clocks maintain information hierarchically.
  - Nodes store local times of a thread + metadata.
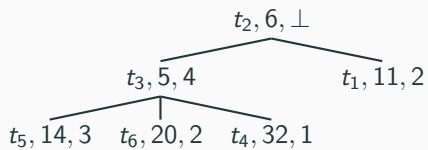  - Tree structure records how information has been obtained transitively.

$$\mathbb{C}_{t_4} = [\overset{t_1}{1}, \overset{t_2}{2}, \overset{t_3}{13}, \overset{t_4}{9}] \qquad \longrightarrow \qquad$$

$$t_4, 9, \bot$$
$$t_3, 13, 5 \quad t_2, 2, 1$$
$$t_1, 1, 1$$

- Only slightly more information is stored compared to vector clocks.

$$\downarrow$$

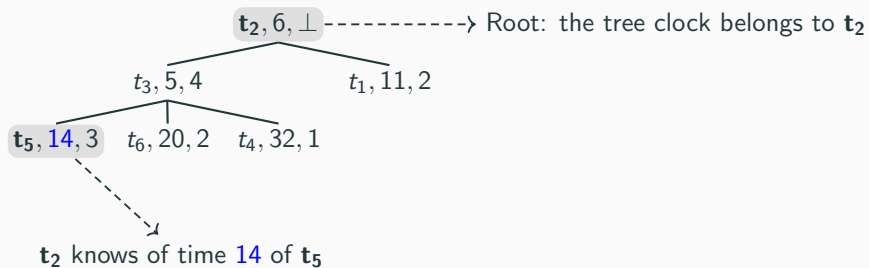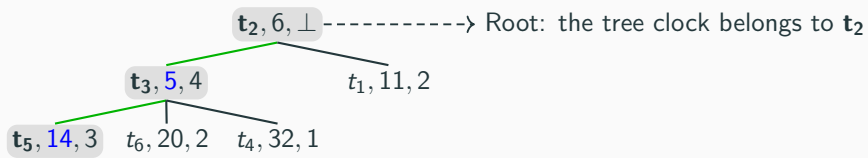**Enough to enable tree clocks to support sub-linear join and (monotone) copy!**

## Tree Clock Data Structure



$t_2, 6, \bot$

$t_3, 5, 4$  $t_1, 11, 2$

$t_5, 14, 3$  $t_6, 20, 2$  $t_4, 32, 1$

$t_2, 6, \perp$ ----------> Root: the tree clock belongs to $t_2$

$t_3, 5, 4$      $t_1, 11, 2$

$t_5, 14, 3$    $t_6, 20, 2$    $t_4, 32, 1$

$t_2, 6, \perp$ - - - - - - - - - $\rightarrow$ Root: the tree clock belongs to $t_2$

$t_3, 5, 4$　　　　$t_1, 11, 2$

$t_5, 14, 3$　$t_6, 20, 2$　$t_4, 32, 1$

$t_2$ knows of time $14$ of $t_5$

$t_2, 6, \bot$ ------------> Root: the tree clock belongs to $t_2$

$t_3, 5, 4$      $t_1, 11, 2$

$t_5, 14, 3$    $t_6, 20, 2$    $t_4, 32, 1$

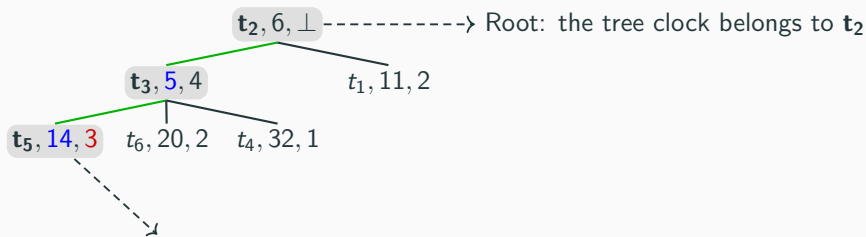$t_2$ knows of time 14 of $t_5$
  $\hookrightarrow$ It learned this transitively, by learning of time 5 of $t_3$
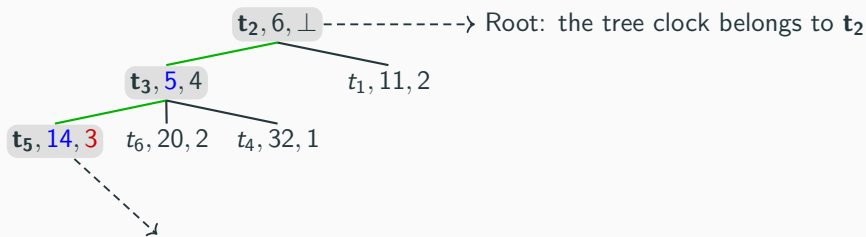
# Tree Clock Data Structure



$\mathbf{t_2}, 6, \bot$ ----------> Root: the tree clock belongs to $\mathbf{t_2}$

$\mathbf{t_3}, 5, 4$      $t_1, 11, 2$

$\mathbf{t_5}, 14, 3$    $t_6, 20, 2$    $t_4, 32, 1$

$\mathbf{t_2}$ knows of time 14 of $\mathbf{t_5}$

↪ It learned this transitively, by learning of time 5 of $\mathbf{t_3}$

↪ $\mathbf{t_3}$ learned of time 14 of $\mathbf{t_5}$ when $\mathbf{t_3}$'s time was 3

# Tree Clock Data Structure



$t_2, 6, \perp$ ----------→ Root: the tree clock belongs to $t_2$

$t_3, 5, 4$          $t_1, 11, 2$

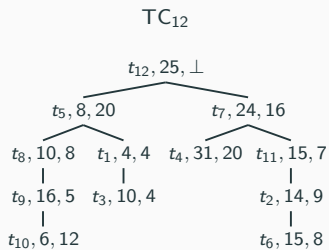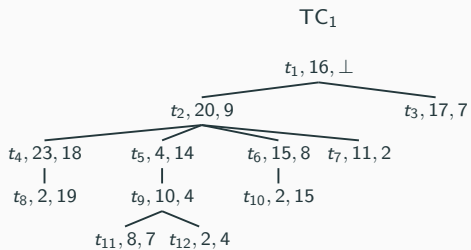$t_5, 14, 3$   $t_6, 20, 2$   $t_4, 32, 1$

$t_2$ knows of time $14$ of $t_5$
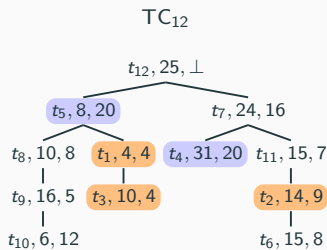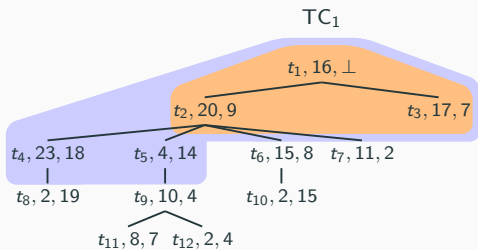  ↪ It learned this transitively, by learning of time $5$ of $t_3$
  ↪ $t_3$ learned of time $14$ of $t_5$ when $t_3$'s time was $3$

**This structure allows for sub-linear time join and (monotone) copy.**

TC$_1$

$t_1, 16, \perp$
$t_2, 20, 9$          $t_3, 17, 7$
$t_4, 23, 18$    $t_5, 4, 14$    $t_6, 15, 8$   $t_7, 11, 2$
$t_8, 2, 19$     $t_9, 10, 4$    $t_{10}, 2, 15$
$t_{11}, 8, 7$   $t_{12}, 2, 4$

TC$_{12}$

$t_{12}, 25, \perp$
$t_5, 8, 20$          $t_7, 24, 16$
$t_8, 10, 8$   $t_1, 4, 4$   $t_4, 31, 20$   $t_{11}, 15, 7$
$t_9, 16, 5$   $t_3, 10, 4$                    $t_2, 14, 9$
$t_{10}, 6, 12$                                 $t_6, 15, 8$

$TC_1$

$t_1, 16, \bot$

$t_2, 20, 9$          $t_3, 17, 7$

$t_4, 23, 18$    $t_5, 4, 14$       $t_6, 15, 8$    $t_7, 11, 2$

$t_8, 2, 19$      $t_9, 10, 4$       $t_{10}, 2, 15$

$t_{11}, 8, 7$    $t_{12}, 2, 4$

$TC_{12}$

$t_{12}, 25, \bot$

$t_5, 8, 20$          $t_7, 24, 16$

$t_8, 10, 8$   $t_1, 4, 4$    $t_4, 31, 20$   $t_{11}, 15, 7$

$t_9, 16, 5$   $t_3, 10, 4$                      $t_2, 14, 9$

$t_{10}, 6, 12$                                  $t_6, 15, 8$

■ Accessed + Updated
■ Accessed

$TC_{12}.Join(TC_1)$

$t_{12}, 25, \bot$

$t_1, 16, 25$     $t_5, 8, 20$      $t_7, 24, 16$

$t_2, 20, 9$  $t_3, 17, 7$  $t_8, 10, 8$   $t_4, 31, 20$  $t_{11}, 15, 7$

$t_6, 15, 8$       $t_9, 16, 5$

$t_{10}, 6, 12$

11

$TC_1$

$TC_{12}$

$t_1, 16, \bot$

$t_2, 20, 9$　　　$t_3, 17, 7$

$t_4, 23, 18$　$t_5, 4, 14$　　$t_6, 15, 8$　$t_7, 11, 2$

$t_8, 2, 19$　　$t_9, 10, 4$　　$t_{10}, 2, 15$

$t_{11}, 8, 7$　$t_{12}, 2, 4$

$t_{12}, 25, \bot$

$t_5, 8, 20$　　　　$t_7, 24, 16$

$t_8, 10, 8$　$t_1, 4, 4$　$t_4, 31, 20$　$t_{11}, 15, 7$

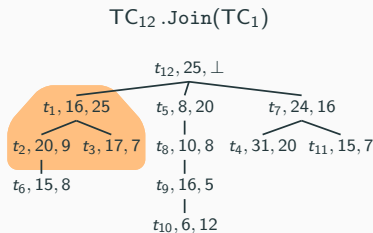$t_9, 16, 5$　$t_3, 10, 4$　　　　　$t_2, 14, 9$

$t_{10}, 6, 12$　　　　　　　　$t_6, 15, 8$

■ Accessed + Updated
■ Accessed

**Performed join without accessing the whole tree!**

$TC_{12}.\mathrm{Join}(TC_1)$

$t_{12}, 25, \bot$

$t_1, 16, 25$　$t_5, 8, 20$　　$t_7, 24, 16$

$t_2, 20, 9$　$t_3, 17, 7$　$t_8, 10, 8$　$t_4, 31, 20$　$t_{11}, 15, 7$

$t_6, 15, 8$　　　　$t_9, 16, 5$

$t_{10}, 6, 12$

11

## $TC_1$

$t_1, 16, \bot$

$t_2, 20, 9$     $t_3, 17, 7$

$t_4, 23, 18$    $t_5, 4, 14$     $t_6, 15, 8$   $t_7, 11, 2$

$t_8, 2, 19$     $t_9, 10, 4$     $t_{10}, 2, 15$

$t_{11}, 8, 7$   $t_{12}, 2, 4$

## $TC_{12}$

$t_{12}, 25, \bot$

$t_5, 8, 20$     $t_7, 24, 16$

$t_8, 10, 8$   $t_1, 4, 4$    $t_4, 31, 20$   $t_{11}, 15, 7$

$t_9, 16, 5$   $t_3, 10, 4$        $t_2, 14, 9$

$t_{10}, 6, 12$            $t_6, 15, 8$

■ Accessed + Updated
■ Accessed

**Performed join without accessing the whole tree!**

### $TC_{12}.\text{Join}(TC_1)$

$t_{12}, 25, \bot$

$t_1, 16, 25$    $t_5, 8, 20$     $t_7, 24, 16$

$t_2, 20, 9$   $t_3, 17, 7$   $t_8, 10, 8$   $t_4, 31, 20$   $t_{11}, 15, 7$

$t_6, 15, 8$        $t_9, 16, 5$

$t_{10}, 6, 12$

|  | $t_1$ | $t_2$ | $t_3$ | $t_4$ | $t_5$ | $t_6$ | $t_7$ | $t_8$ | $t_9$ | $t_{10}$ | $t_{11}$ | $t_{12}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $\mathbb{C}_{t_1} =[$ | 16 | 20 | 17 | 23 | 4 | 15 | 11 | 8 | 10 | 2 | 8 | 7 ] |

### Drop-in Replacement

**Algorithm:** Happens-Before with **Vector Clocks**.

1 **procedure** acquire($t$, $\ell$)
2 $\quad$ $\mathbb{C}_t$.VectorClockJoin($\mathbb{C}_\ell$)

3 **procedure** release($t$, $\ell$)
4 $\quad$ $\mathbb{C}_\ell$.VectorClockCopy($\mathbb{C}_t$)

---

**Algorithm:** Happens-Before with **Tree Clocks**.

1 **procedure** acquire($t$, $\ell$)
2 $\quad$ $\mathbb{C}_t$.TreeClockJoin($\mathbb{C}_\ell$)

3 **procedure** release($t$, $\ell$)
4 $\quad$ $\mathbb{C}_\ell$.TreeClockCopy($\mathbb{C}_t$)

What is the optimal data structure for Happens-Before?

## Tree Clock Optimality

What is the optimal data structure for Happens-Before?

**Treeclock optimality for Happens-Before**

No other data structure can offer asymptotically better performance.

## Tree Clock Optimality

What is the optimal data structure for Happens-Before?

**Treeclock optimality for Happens-Before**

No other data structure can offer asymptotically better performance.

- Tree clocks perform at most 3 times more work than necessary.

$TC_1$

Accessed + Updated
Accessed

$t_1, 16, \perp$

$t_2, 20, 9$      $t_3, 17, 7$

$t_4, 23, 18$    $t_5, 4, 14$    $t_6, 15, 8$   $t_7, 11, 2$

$t_8, 2, 19$    $t_9, 10, 4$    $t_{10}, 2, 15$

$t_{11}, 8, 7$   $t_{12}, 2, 4$

**Vector time work** $\mathsf{VTWork}(\sigma)$

$\mathsf{VTWork}(\sigma) =$ the *smallest* number of data-structure accesses for processing $\sigma$

# Tree Clock Optimality



Accessed + Updated
Accessed

$TC_1$

$t_1, 16, \bot$

$t_2, 20, 9$      $t_3, 17, 7$

$t_4, 23, 18$    $t_5, 4, 14$    $t_6, 15, 8$   $t_7, 11, 2$

$t_8, 2, 19$    $t_9, 10, 4$    $t_{10}, 2, 15$

$t_{11}, 8, 7$   $t_{12}, 2, 4$

**Vector time work** $\mathsf{VTWork}(\sigma)$

$\mathsf{VTWork}(\sigma) =$ the *smallest* number of data-structure accesses for processing $\sigma$

**Tree clock work** $\mathsf{TCWork}(\sigma)$

$\mathsf{TCWork}(\sigma) =$ the total number of *tree* clock entries accessed for processing $\sigma$

14

# Tree Clock Optimality



$\mathsf{TC}_1$

Accessed + Updated
Accessed

$t_1, 16, \bot$
$t_2, 20, 9$
$t_3, 17, 7$
$t_4, 23, 18$
$t_5, 4, 14$
$t_6, 15, 8$
$t_7, 11, 2$
$t_8, 2, 19$
$t_9, 10, 4$
$t_{10}, 2, 15$
$t_{11}, 8, 7$
$t_{12}, 2, 4$

**Vector time work** $\mathsf{VTWork}(\sigma)$

$\mathsf{VTWork}(\sigma) =$ the *smallest* number of data-structure accesses for processing $\sigma$

**Tree clock work** $\mathsf{TCWork}(\sigma)$

$\mathsf{TCWork}(\sigma) =$ the total number of *tree* clock entries accessed for processing $\sigma$

**Vector clock work** $\mathsf{VCWork}(\sigma)$

$\mathsf{VCWork}(\sigma) =$ the total number of *vector* clock entries accessed for processing $\sigma$

# Data Structure Optimality for Happens-Before



$$\mathsf{VCWork}(\sigma) \leq \mathcal{T} \cdot \mathsf{VTWork}(\sigma)$$

$\mathsf{VCWork}(\sigma)$ can be $\mathcal{T}$ times worse

# Data Structure Optimality for Happens-Before



$\text{VCWork}(\sigma) \leq \mathcal{T} \cdot \text{VTWork}(\sigma)$

$\text{VCWork}(\sigma)$ can be $\mathcal{T}$ times worse

**Theorem**

$\text{TCWork}(\sigma) \leq 3 \cdot \text{VTWork}(\sigma)$

Tree clocks are (asymptotically) VT-optimal!

## Beyond Happens-Before

- Tree clocks can be used to compute other partial orders.

## Beyond Happens-Before

- Tree clocks can be used to compute other partial orders.

Schedulable-Happens-Before (SHB)

- Used in sound data race detection[1]



[1]U. Mathur, D. Kini, M. Viswanathan. What Happens-after the First Race? Enhancing the Predictive Power of Happens-before Based Dynamic Race Detection. OOPSLA'18.

## Beyond Happens-Before

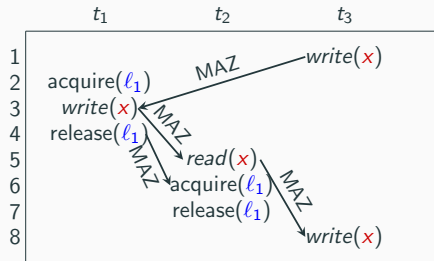- Tree clocks can be used to compute other partial orders.

Schedulable-Happens-Before (SHB)

- Used in sound data race detection[1]

Mazurkiewicz (MAZ)

- Used in dynamic partial order reduction in model checking of concurrent programs[2].





---

[1] U. Mathur, D. Kini, M. Viswanathan. What Happens-after the First Race? Enhancing the Predictive Power of Happens-before Based Dynamic Race Detection. OOPSLA'18.

[2] C. Flanagan, P. Godefroid. Dynamic Partial-Order Reduction for Model Checking Software. POPL'05.
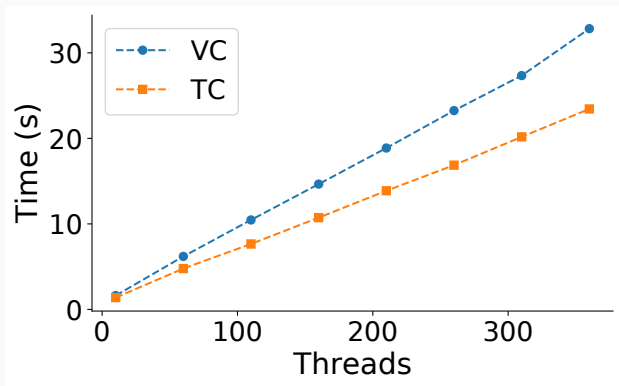
## Experimental Results



- 153 benchmark traces
  - Based on standard Java and OpenMP benchmark suites.
- Implemented HB, SHB and MAZ with both tree clocks and vector clocks.
- Measured the time in the following tasks:
  - Compute the partial order.
  - Perform the race-detection analysis.

## Experimental Results

- 153 benchmark traces
  - Based on standard Java and OpenMP benchmark suites.
- Implemented HB, SHB and MAZ with both tree clocks and vector clocks.
- Measured the time in the following tasks:
  - Compute the partial order.
  - Perform the race-detection analysis.

|  | Mazurkiewicz | Schedulable-Happens-Before | Happens-Before |
|---|---|---|---|
| Partial Order | 2.02× | 2.66× | 2.97× |
| Partial Order + Analysis | 1.49× | 1.80× | 1.11× |

Significant speedup by just replacing vector clocks with tree clocks!
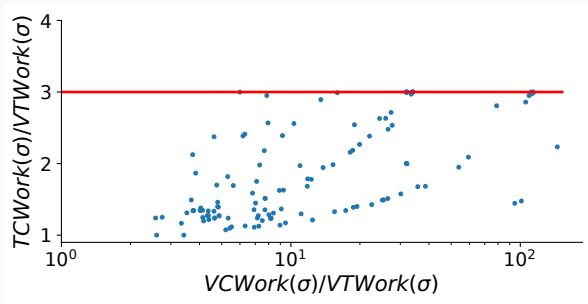
## Experimental Results - Scalability

- Controlled experiment: threads randomly communicate over a single lock.
  - Theoretical speedup: $4\times$
  - Observed speedup: $1.33\times$

**Theorem**

$\mathsf{TCWork}(\sigma) \leq 3 \cdot \mathsf{VTWork}(\sigma)$

## Conclusion

<u>Tree clocks</u>:

- Drop-in replacement of vector clocks.

## Conclusion

Tree clocks:

- Drop-in replacement of vector clocks.
- Tree clocks support join and (monotone) copy in sub-linear time.

## Conclusion

Tree clocks:

- Drop-in replacement of vector clocks.
- Tree clocks support join and (monotone) copy in sub-linear time.
- Optimal data structure for Happens-Before (asymptotically).

## Conclusion

Tree clocks:

- Drop-in replacement of vector clocks.
- Tree clocks support join and (monotone) copy in sub-linear time.
- Optimal data structure for Happens-Before (asymptotically).
- Experimental results confirm the potential of tree clocks.

## Conclusion

Tree clocks:

- Drop-in replacement of vector clocks.
- Tree clocks support join and (monotone) copy in sub-linear time.
- Optimal data structure for Happens-Before (asymptotically).
- Experimental results confirm the potential of tree clocks.

Thank you!