

*Plausible Clocks with Bounded Inaccuracy*

---

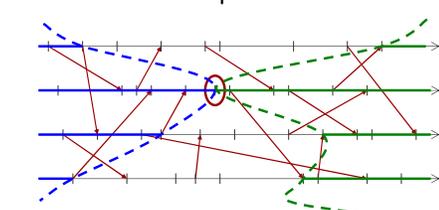
**Brad Moore, Paul Sivilotti**

Computer Science & Engineering  
The Ohio State University  
paolo@cse.ohio-state.edu



## Causality & Concurrency

- $a \rightarrow b \Leftrightarrow$  exists a path from a to b



- $a || b \Leftrightarrow \neg(a \rightarrow b) \wedge \neg(b \rightarrow a)$

Paul A.G. Sivilotti, DISC 2005

## Time-Stamping Systems

- A TSS is an algorithm for assigning
  - "time stamps" to events,  $T(a)$
  - "time tags" to messages
- Implementing a TSS requires
  - STAMP()
    - Input: previous local stamp, incoming tag
    - Output: new local stamp
  - TAG()
    - Input: previous local stamp
    - Output: tag to add to outgoing message
  - COMP()
    - Input: two time stamps
    - Output: one of { <, =, || }



Paul A.G. Sivilotti, DISC 2005

## Example TSS: Lamport Clocks

- Time stamps are integers
- Time tags are integers
- STAMP:
  - for local and send events, increase stamp
  - for receive events, use max of stamp & tag
- TAG:
  - use local time stamp to tag messages
- COMP:
  - < and =, integer comparison
  - || when same value but on different processes



Paul A.G. Sivilotti, DISC 2005

## Example TSS: Vector Clocks

- Time stamps are arrays of N integers
- Time tags are arrays of N integers
- STAMP:
  - for local and send events, increase local element
  - for receive events, use element-wise max of stamp & tag
- TAG:
  - use local time stamp to tag messages
- COMP:
  - < when all elements are <= and at least one is <
  - = when all elements are =
  - || when some element is < and another is >



Paul A.G. Sivilotti, DISC 2005

## Plausibility

- Weak clock condition (WCC)
  - $a \rightarrow b \Rightarrow T(a) < T(b)$
- A TSS satisfying the WCC is **plausible** iff it is consistent with underlying hb partial order
  - every pair ordered by TSS is causally related
  - no distinct events are given equal stamps
- Example: Lamport clocks
  - orders (almost) all pairs!
    - exception: same value on different processes
  - every pair identified as concurrent, is indeed concurrent in hb partial order
  - stamps from distinct events can always be distinguished (add process id to stamp)



Paul A.G. Sivilotti, DISC 2005

## Ordering from Lamport Clocks

7

## Application of Plausibility

- Correctly identifying causally related pairs generally necessary for **safety**
  - arbitration (resource allocation, mutex)
  - consistent serialization (cache coherence)
  - playback for distributed debugging
- Correctly identifying concurrent pairs generally important for **performance**
  - consistent cuts require concurrent sets
  - cache consistency and invalidation protocols
  - snapshots (deadlock & termination detection, checkpoints for rollback)

8

## Inaccuracy of a Plausible TSS

- Two kinds of errors:
  - "false concurrency": related events stamped concurrent
    - WCC means a plausible TSS makes of none of these errors
  - "false ordering": concurrent events stamped as ordered
    - plausible clocks may make this error
- "Inaccuracy" measures how many such mistakes (false orderings) are made in a given run
- Fundamental result:
  - perfect accuracy requires  $O(N)$  message overhead
  - (so Vector clocks are optimal)
  - does not scale to large systems
- Motivates natural research question:
  - Can we get good expected-case accuracy with less message overhead?

9

## Designing a Plausible TSS

- Previous approaches:
  - Fix time tag size (ie message overhead)
    - constant size
  - Measure resulting accuracy through simulation
    - expected case analysis

10

## Inversion of Design

- Our approach:
  - Fix inaccuracy
    - some constant upper bound
  - Measure resulting message complexity through simulation
    - expected case analysis

11

## Towards a New Metric: Error Count ( $\delta$ )

- *Error count* for an event is the number of false (pre) orderings, for that event

$$\delta(P, H, b) = \left| \left\{ a \in H :: a \parallel b \wedge a \xrightarrow{P} b \right\} \right|$$

- i.e., number of *blue* dots

12

## Towards a New Metric: Inaccuracy ( $\rho$ )

- Inaccuracy** is the ratio of false orderings, averaged over entire run

$$\rho(P, H) = \frac{2 \times (\sum_{b \in H :: \delta(P, H, b)})}{|\{a, b \in H :: a \parallel b\}|}$$

- i.e., ratio of dots to concurrent events

Paul A.G. Sivillotti, DISC 2005 13

## Introducing a New Metric: Imprecision ( $\psi$ )

- Imprecision** is the max number of false (pre) orderings for a given time stamp

$$\psi(P, s) = (\text{Max } H \in \mathbf{H}(P, s), a \in H : P(a) = s : \delta(P, H, a))$$

- i.e., **worst-case** number of **blue dots**
- independent of particular run (property of a **time stamp**)

Paul A.G. Sivillotti, DISC 2005 14

## Introducing a New Metric: Imprecision ( $\psi$ )

- Imprecision** is the max number of false (pre) orderings for a given time stamp

$$\psi(P, s) = (\text{Max } H \in \mathbf{H}(P, s), a \in H : P(a) = s : \delta(P, H, a))$$

- i.e., **worst-case** number of **blue dots**
- independent of particular run (property of a **time stamp**)

Paul A.G. Sivillotti, DISC 2005 15

## Algorithmic Approach

- Recall Vector Clocks**
  - stamp is an array of values,  $V[1..N]$
  - $V[i]$  is "most recent event" on process  $i$

Paul A.G. Sivillotti, DISC 2005 16

## Compressing Time Tags

- If multiple entries have the same value, they can share an entry in vector
  - problem: unlikely to have exact equality
- Use a **range** for vector entry
  - most recent event guaranteed to be in range
  - the larger the range, the greater the imprecision
- Allow ranges to **grow/shrink** (and number of vector entries to increase/decrease)
  - maintain constant precision
  - message overhead may vary over the run

Paul A.G. Sivillotti, DISC 2005 17

## Three Claims

- The algorithm is a correct plausible TSS
  - satisfies the **WCC** (and distinguishes events)
- The imprecision can be controlled
  - different sources of imprecision: local stamps and tagged messages
  - their combination does not lose too much information
- The algorithm achieves good performance
  - measured in terms of **tag size**
  - expected case** evaluation

Paul A.G. Sivillotti, DISC 2005 18

## Outline of the Talk

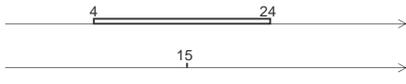
- Background
  - model, causality ( $\rightarrow$ ), concurrency, TSS
- Plausible TSS
  - WCC, accuracy vs precision
- Algorithmic Approach
  - use ranges instead of precise values
- Algorithm Description
  - time intervals
  - implementation of STAMP, TAG, COMP
- Justification of claims
  - correctness (algorithm is a plausible TSS)
  - imprecision is bounded
  - performance is reasonable



19  
Paul A.G. Sivikotti, DISC 2005

## Time Intervals

- Interval given by beginning and end points
  - if these are the same, interval is "precise"



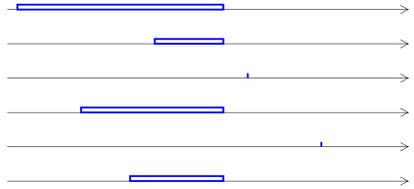
- An interval is associated with each process
  - goal: the "local time" of the most recent causally related event is within the interval
  - interval associated with *this* process is precise
  - local storage is linear in N



20  
Paul A.G. Sivikotti, DISC 2005

## Time Stamps

- A time stamp is a vector of intervals
- Satisfies two invariants:
  - imprecise intervals share the **same end point**
  - precise intervals are all **greater** than this value




21  
Paul A.G. Sivikotti, DISC 2005

## Time Tags

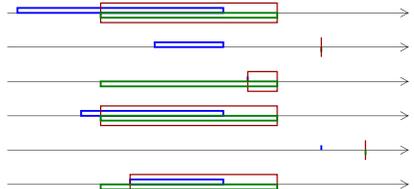
- Subset of time stamps
  - a vector of intervals, with properties A & B
- Satisfies an additional invariant:
  - imprecise intervals share the **same begin point**
  - all imprecise intervals are the same (a "bucket")




22  
Paul A.G. Sivikotti, DISC 2005

## Implementation: STAMP

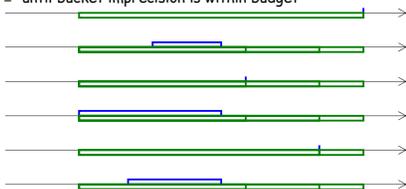
- For local/send events: increase local (precise) interval
- For receive events (stamp + tag)
  - each interval in new stamp defined by:
    - begin = max of begin points
    - end = max of end points
  - increase local (precise) interval




23  
Paul A.G. Sivikotti, DISC 2005

## Implementation: TAG

- Given a time stamp and an imprecision budget:
  - build smallest acceptable tag (ie use as big a bucket as possible)
- Algorithm:
  - start with everything in bucket
  - repeat: move most recent interval from bucket to its own entry
  - until bucket imprecision is within budget




24  
Paul A.G. Sivikotti, DISC 2005

## Implementation: COMP

- Comparing individual intervals  $(i, j)$ 
  - $i = j$ : same **precise** intervals
 
  - $i < j$ : no overlap, end of  $i <$  beg of  $j$ 

  - $i \approx j$ : some overlap
 
  - $i < \approx j$ :  $\neg(j < i)$ 


Paul A.G. Sivillotti, DISC 2005 25

## Implementation: COMP

- Comparing vectors of intervals  $(s, t)$ 
  - $s < t$ 
    - all  $i$  in  $s, j$  in  $t :: i < \approx j$
    - exists an  $i$  in  $s, j$  in  $t :: i < j$
  - $s = t$ 
    - all  $i$  in  $s, j$  in  $t :: i \approx j$
  - $s || t$ 
    - exists an  $i$  in  $s, j$  in  $t :: i < j$
    - exists an  $i$  in  $s, j$  in  $t :: j < i$

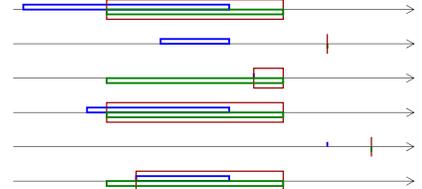
Paul A.G. Sivillotti, DISC 2005 26

## Justifying the Claims

- The algorithm is a correct plausible TSS
  - stamp and tag invariants [A-C]
  - $a \rightarrow b \Rightarrow T(a) < T(b)$  [WCC]
  - $a = b \Leftrightarrow T(a) = T(b)$  [no 2 equal stamps]
- The imprecision is bounded (worst case)
- The algorithm achieves good performance

Paul A.G. Sivillotti, DISC 2005 27

## Claim 1: Correctness

- Part (i): stamp and tag invariants
  - imprecise intervals share an end point
  - precise intervals are all greater than this value
- Proof intuition
  - new shared end point = max of old shared end points

Paul A.G. Sivillotti, DISC 2005 28

## Claim 1: Correctness

- Part (ii): WCC holds
  - $a \rightarrow b \Rightarrow T(a) < T(b)$
- Surprising complication:  $\approx$  is not transitive
  - $i < \approx j$
  - $j < \approx k$
  - $\neg(i < \approx k)$
- Proof intuition
  - along every **actual chain**, time stamps are non decreasing
  - every happens-before pair is joined by a chain

Paul A.G. Sivillotti, DISC 2005 29

## Claim 1: Correctness

- Part (iii): No two equal stamps
  - $a = b \Leftrightarrow T(a) = T(b)$
- Surprising complication:
  - stamps are equal when all entries overlap
- Proof intuition:
  - local interval is always precise
  - for events on same process, the local interval is unique
  - for events on different processes, both local intervals can not be identical
    - equality of local interval implies causality
    - hence, equality of both implies a cycle in causality

Paul A.G. Sivillotti, DISC 2005 30

## Claim 2: Guaranteed Bound

- Surprising since:
  - stamp intervals can **increase** on receive!
  - tagging **increases** imprecision
- Key observation:
  - after receive, each stamp interval smaller than **either old stamp or tag**
- Important properties:
  - $\psi(\text{stamp}) \leq \max(\psi(\text{old stamp}), \psi(\text{tag}))$
  - $\psi(\text{tag}) \leq \text{BOUND}$
- Result:
  - stable.  $\psi(\text{stamp}) \leq \text{BOUND}$



31

Paul A.G. Sivilotti, DISC 2005

## Claim 3: Good Performance

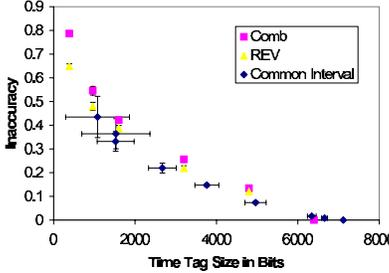
- Comparison challenges
  - trade-off is similar (message size vs errors)
  - but imprecision is **not** accuracy!
- Simulation results
  - varying number of processes
  - fixed topology
    - some high degree nodes some low
  - varying intercommunication delays



32

Paul A.G. Sivilotti, DISC 2005

## Performance: Good News (Inaccuracy vs. Message Size)



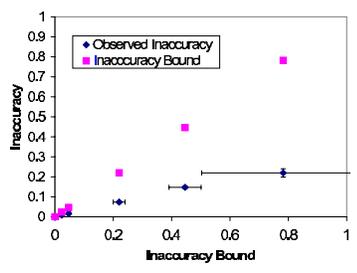
Time Tag Size (Bits)	Comb (Inaccuracy)	REV (Inaccuracy)	Common Interval (Inaccuracy)
0	0.8	0.6	0.4
1000	0.5	0.4	0.35
2000	0.4	0.35	0.3
3000	0.25	0.2	0.2
4000	0.15	0.1	0.15
5000	0.1	0.05	0.1
6000	0.05	0.02	0.05
7000	0.02	0.01	0.02
8000	0.01	0.005	0.01



33

Paul A.G. Sivilotti, DISC 2005

## Performance: Good (?) News (Observed vs Bound)



Inaccuracy Bound	Observed Inaccuracy
0	0.05
0.1	0.08
0.2	0.1
0.3	0.15
0.4	0.15
0.5	0.15
0.6	0.15
0.7	0.15
0.8	0.15
0.8	0.8



34

Paul A.G. Sivilotti, DISC 2005

## Conclusions

- Plausible clocks
  - provide tradeoff between message size and accuracy
- Contribution:** Imprecision metric
  - property of a time stamp (run independent)
  - worst-case behavior
  - unbounded for any constant message size
- Contribution:** Bounded imprecision plausible clocks
  - time stamps have intervals (messages have buckets)
  - message overhead grows and shrinks as necessary
  - guarantee on amount of imprecision (hence inaccuracy)
- Contribution:** Promising performance
  - actual accuracy generally better than imprecision bound



35

Paul A.G. Sivilotti, DISC 2005

## Plausible Clocks with Bounded Inaccuracy

**Brad Moore, Paul Sivilotti**

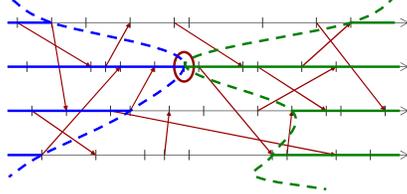
Computer Science & Engineering  
The Ohio State University  
paolo@cse.ohio-state.edu



## Bounding Inaccuracy

■ Concurrency ratio,  $\epsilon$ :  $\epsilon(H) = \frac{1}{2} \times \frac{|a, b \in H :: a \parallel b|}{|H|}$

■ i.e., (avg) number of events in concurrent window



■ For "well-behaved" computations, this is constant



Paul A.G. Sivillotti, DISC 2005

37

## Bounding Inaccuracy (Cont'd)

■ Rewrite inaccuracy using  $\epsilon(H)$

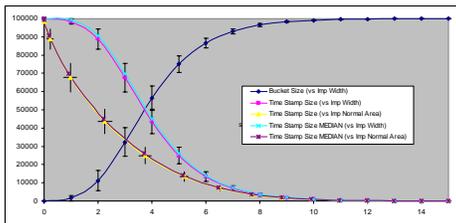
$$\begin{aligned} \rho(P, H) &= \frac{2 \times (\sum b \in H :: \delta(P, H, b))}{|a, b \in H :: a \parallel b|} \\ &\leq \frac{2 \times (\sum b \in H :: w(P, P.stamp(b)))}{|a, b \in H :: a \parallel b|} \\ &= \frac{1}{\epsilon(H)} \times \frac{(\sum b \in H :: w(P, P.stamp(b)))}{|H|} \\ &\leq \frac{1}{\epsilon(H)} \times BOUND \end{aligned}$$



Paul A.G. Sivillotti, DISC 2005

38

## Performance: Bucket Size



Paul A.G. Sivillotti, DISC 2005

39