

CSCI 4210/6210 Parallel & Distributed Simulation

PDES: Time Warp Mechanism
Distributed Snapshots and Fossil Collection



Maria Hybinette, UGA

Global Virtual Time

Wallclock time $T(GVT_t)$ during the execution of a Time Warp simulation is defined as the **minimum time stamp among all unprocessed and partially processed messages and anti-messages in the system at wall-clock T .**

Maria Hybinette, UGA

Outline

- Consistent Cuts
 - » Cut points
 - » Cut messages
 - » Cut values
- Mattern's GVT Algorithm
 - » Colors
 - » Vector counters
 - » Pipelined algorithm
- Fossil Collection
- Thursday: We will do an additional example using Mattern's Algorithms to determine GVT.

Maria Hybinette, UGA

Review: Samadi's Algorithm

- Transient message problem:
 - » **Solution:** Message acknowledgements
- Simultaneous message problem:
 - » **Solution:** Mark acknowledgements sent after reporting local minimum
 - » **Caveat:** Just message acks are not enough (marked message acks are needed).
- Overhead:
 - » Message acknowledgments:
 - Message acknowledgment for
 - each message and
 - anti-message.

Maria Hybinette, UGA

Mattern's Algorithm

- Asynchronous
 - » Executes in background concurrent with time warp execution (does not require the simulation to "freeze" (i.e., block the LPs).
- Avoids message acknowledgements
- Approach: Based on techniques for creating distributed snapshots (consistent cut)
 - » We will see what it means to be a consistent cut

- » Can some asynchronous algorithms compute *exact* GVT(t)?
- » What about synchronous algorithms?

Maria Hybinette, UGA

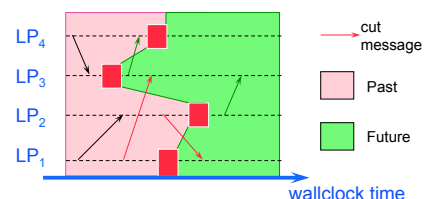
Consistent Cuts

Cut point:

Cut:

Cut message:

Consistent cut:



Cut value:

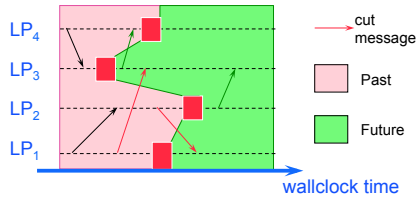
Consistent Cuts

Cut point: an instant dividing computation into **past** and **future**

Cut:

Cut message:

Consistent cut:



Cut value:

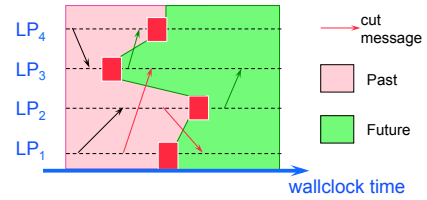
Consistent Cuts

Cut point: an instant dividing computation into **past** and **future**

Cut: set of cut points, one per processor

Cut message:

Consistent cut:



Cut value:

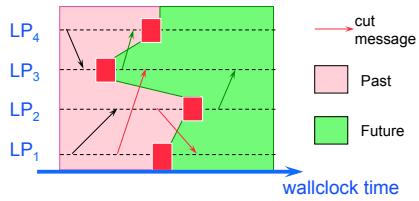
Consistent Cuts

Cut point: an instant dividing computation into **past** and **future**

Cut: set of cut points, one per processor

Cut message: a message that was sent in the past, and received in the future

Consistent cut:



Cut value:

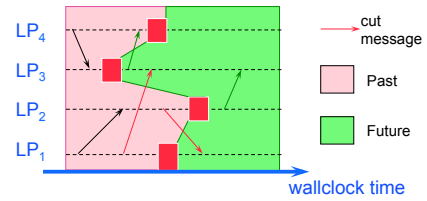
Consistent Cuts

Cut point: an instant dividing computation into **past** and **future**

Cut: set of cut points, one per processor

Cut message: a message that was sent in the past, and received in the future

Consistent cut: a cut where all messages crossing the cut are cut messages



Cut value:

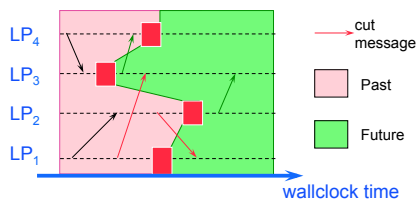
Consistent Cuts

Cut point: an instant dividing computation into **past** and **future**

Cut: set of cut points, one per processor

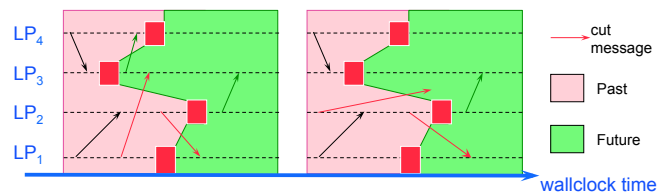
Cut message: a message that was sent in the past, and received in the future

Consistent cut: a cut where all messages crossing the cut are cut messages



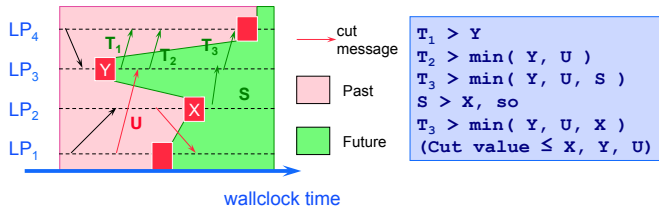
Cut value: minimum among (1) local minimum of each LP at its cut point and (2) time stamp of cut messages

Cuts: Divides Past and Future



- **Consistent Cuts:** Includes local state at its cut-point & all its transient messages.
- **Observation:** Time stamp of a message sent after a cut point at wallclock time T must be at least as large as the minimum of:
 - » the smallest time stamp of any unprocessed event in the processor at time T
 - » the smallest time stamp of any message received by the processor after time T.
- **GVT** must be smaller than or equal to both of these quantities

Observation 1

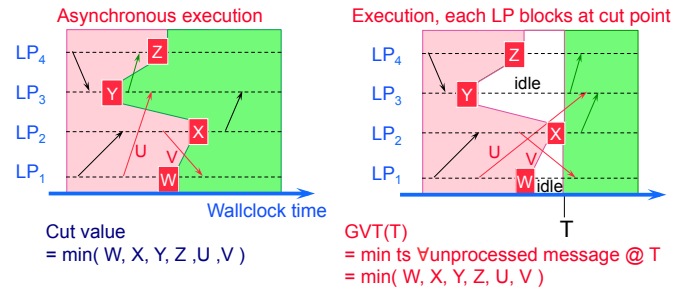


Any message crossing cut from future to past must have a time stamp **greater than the cut value**, so they can be ignored when computing the cut value

Message generated by an LP after its cut point must have time stamp **greater than the minimum of**

- » The LP's local minimum at its cut point
- » The time stamp of messages received after the cut point

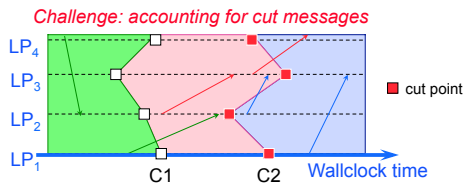
Observation 2



- Cut value equal to GVT(T) using **synchronous GVT** algorithm (freeze LPs: no new computations nor message sends/receives).
 - » Events generated after cut have time stamp $>$ cut value

- **Cut value can be used as a GVT value**

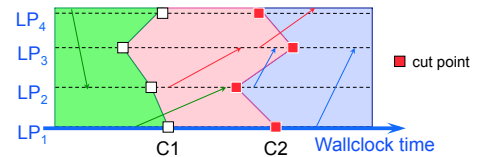
Mattern's GVT Algorithm



Approach:

- Construct two cuts C1, C2, approximate cut value along C2
 - » Organize processes in ring, pass token around ring
- Ensure no message that crosses C1 also cross C2
 - » Color LPs, change LP color at each cut point
 - » Color (green/red) each message to that of LP sending message (message tag)
 - » Maintain send/receive message counters
- GVT = $\min(\text{local min along C2, time stamp of red messages})$

Algorithm Overview



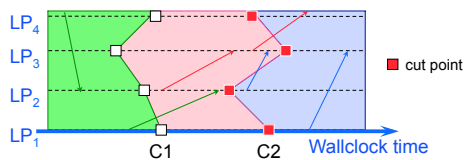
The first cut:

- » Changes color of each process (green to red)
- » Determine number of green messages sent to each process

The second cut:

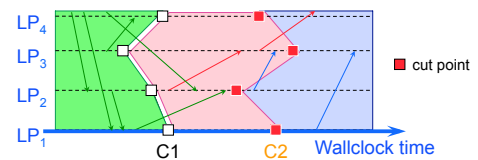
- » Each process makes sure all green messages sent to it have been received before laying down a cut point
- » Compute global minimum (GVT value)

How does an LP know it has received all its green messages?



- LP_i maintains vector $v_i[1:N]$, where $N = \#LPs$
 - » $v_i[i] = \text{number of green messages received by } LP_i$
 - » $v_i[s] = \text{number of green messages sent by } LP_i \text{ to } LP_s$
- C2: LP_i cannot pass token until
 - » $v_i[i] = \sum v_s[i]$ (summed over all $s \neq i$)
- C1: Token includes vector to accumulate send counters

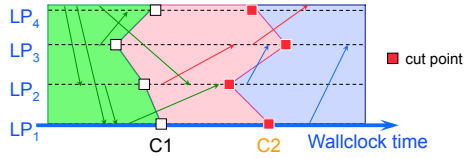
Example: Vector Counters



Vector counters for green messages (at C2) $i = j$ received:

	v_1	v_2	v_3	v_4
$v_1[4] = 0$				

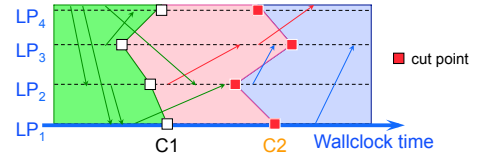
Example: Vector Counters



Vector counters for green messages (at C2) $i = j$ received:

v_1	v_2	v_3	v_4
$v_1[4] = 0$			
$v_1[3] = 0$			

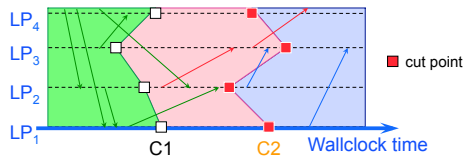
Example: Vector Counters



Vector counters for green messages (at C2) $i = j$ received:

v_1	v_2	v_3	v_4
$v_1[4] = 0$			
$v_1[3] = 0$			
$v_1[2] = 1$			
.			

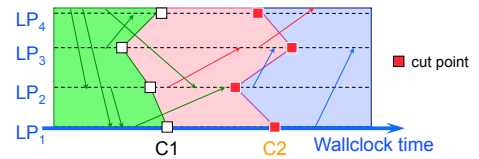
Example: Vector Counters



Vector counters for green messages (at C2) $i = j$ received:

v_1	v_2	v_3	v_4
$v_1[4] = 0$			
$v_1[3] = 0$			
$v_1[2] = 1$			
$v_1[1] = -2$			

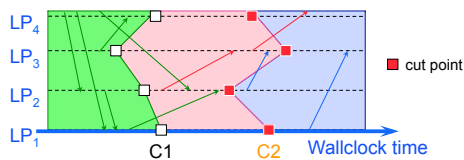
Example: Vector Counters



Vector counters for green messages (at C2) $i = j$ received:

v_1	v_2	v_3	v_4
$v_1[4] = 0$	$v_2[4] = 0$		
$v_1[3] = 0$	$v_2[3] = 0$		
$v_1[2] = 1$	$v_2[2] = -3$		
$v_1[1] = -2$	$v_2[1] = 1$		

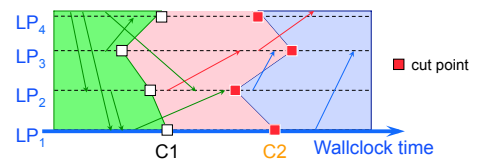
Example: Vector Counters



Vector counters for green messages (at C2) $i = j$ received:

v_1	v_2	v_3	v_4
$v_1[4] = 0$	$v_2[4] = 0$	$v_3[4] = 1$	
$v_1[3] = 0$	$v_2[3] = 0$	$v_3[3] = 0$	
$v_1[2] = 1$	$v_2[2] = 3$	$v_3[2] = 0$	
$v_1[1] = -2$	$v_2[1] = 1$	$v_3[1] = 0$	

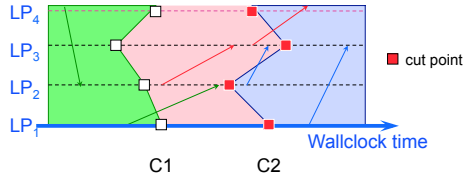
Example: Vector Counters



Vector counters for green messages (at C2) $i = j$ received:

v_1	v_2	v_3	v_4
$v_1[4] = 0$	$v_2[4] = 0$	$v_3[4] = 1$	$v_4[4] = -1$
$v_1[3] = 0$	$v_2[3] = 0$	$v_3[3] = 0$	$v_4[3] = 0$
$v_1[2] = 1$	$v_2[2] = -3$	$v_3[2] = 0$	$v_4[2] = 2$
$v_1[1] = -2$	$v_2[1] = 1$	$v_3[1] = 0$	$v_4[1] = 1$

Mattern's GVT Algorithm



- **Local Variables (in each logical process LP_i):**
 - » T_{red} = min time stamp among red messages sent by LP (even non-cut red messages!)
 - » $V_i[1:N]$ = message send / receive counters
- **Token: CMsg**
 - » $CMsg_T_{min}$ = accumulator, smallest local minimum so far
 - » $CMsg_T_{red}$ = accumulator, smallest red message time stamp so far
 - » $CMsg_Count[1:N]$ = # messages sent to each LP

Mattern's GVT Algorithm

- Message send by green logical process from LP_i to LP_j

$$V_i[j] = V_i[j] + 1$$
- LP_i receives a green message

$$V_i[i] = V_i[i] - 1$$
- Control message, first cut:
 - Change color of process to red
 - $CMsg_Count = CMsg_Count + V_i$
 - Forward control message to next process in ring
- Message send with time stamp ts by a red LP

$$T_{red} = \min(T_{red}, ts)$$
- Control message, second cut:
 - wait until $V_i[i] = CMsg_Count[i]$ i.e., #received = #sent
 - $CMsg_T_{min} = \min(CMsg_T_{min}, T_{min})$
 - $CMsg_T_{red} = \min(CMsg_T_{red}, T_{red})$
 - forward token to next process in ring

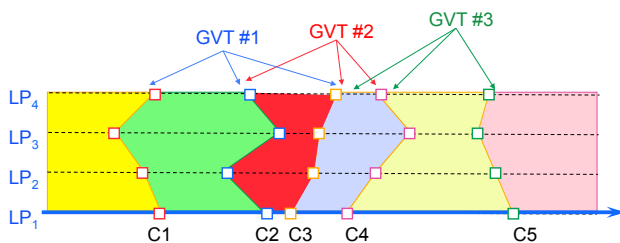
Fossil Collection

- **Batch fossil collection**
 - » After GVT computation, scan through list of LPs mapped to processor to reclaim memory and commit I/O operations
 - » May be time consuming if many LPs
- **On-the-fly Fossil Collection**
 - » After processing event, place memory into "free memory" list
 - » Before allocating memory, check that time stamp is less than GVT before reusing memory

Summary

- Consistent cuts
- Cut value can be used as an estimate of GVT
 - » Local minimum at each LP
 - » Cut messages
- Construct second consistent cut
 - » Coloring LPs, messages
 - » Vector counter to determine when an LP has received all relevant cut messages
- Pipeline GVT computation, continuously circulating token
- Numerous variations
 - » Could implement cuts with other communication topologies, e.g., butterfly
 - » Other ways to deal with transient messages, e.g., global count and abort/retry mechanism for second cut, etc.

Distributing GVT Values & Pipelining



- **Pipelined execution**
 - Overlap successive GVT computations: first GVT uses C1, C2, C3, second uses C2, C3, C4, etc.
 - Each cut computes a new GVT value
 - Continuously circulate GVT token