

Lazy Snapshots: A Taxonomy and Performance Study

Zhijun Liu and Paul Sivilotti*
(and Nigamanth Sridhar)

Computer Science and Engineering
The Ohio State University

*currently visiting faculty at University of Minnesota



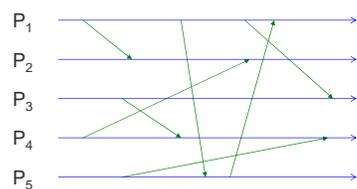
Background: Checkpointing

- Checkpoint = recording state of system
- Utility:
 - Fault recovery
 - Debugging
 - Stable property detection
- State of a distributed system consists of:
 - local state of individual processes
 - state of channels (ie messages in transit)
- Simple algorithm:
 - Every process records its local state at time t
 - Channel state?
- Challenge:
 - No shared clock



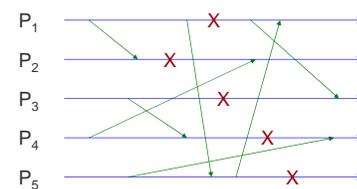
Paul A.G. Sivilotti, PDCS 2005

Consistent Checkpoints



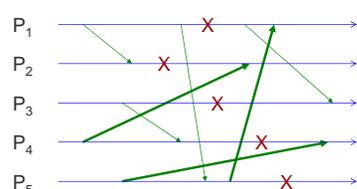

Paul A.G. Sivilotti, PDCS 2005

Consistent Checkpoints: Local State



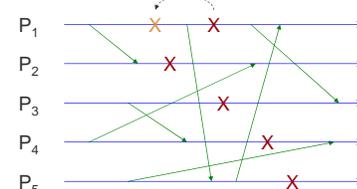

Paul A.G. Sivilotti, PDCS 2005

Consistent Checkpoints: Channel State




Paul A.G. Sivilotti, PDCS 2005

Consistent Checkpoints: Requirement



- The cut must be "input-closed":
Every message recorded as received,
must also have been recorded as sent



Paul A.G. Sivilotti, PDCS 2005

Consistency:
No *after* message should be received *before* RLS

Paul A.G. Sivillotti, PDCS 2005

Completeness (channel state):
before messages rec'ed *after* RLS are in transit

Paul A.G. Sivillotti, PDCS 2005

A Cooperative Snapshot Alg

- Use **marker** to separate *before* messages from *after* messages
- Assumption: FIFO channels

Paul A.G. Sivillotti, PDCS 2005

Rule 1:
Record local state *immediately* when *first* marker arrives

Rule 2:
Record channel state *after* RLS and *until* marker arrives

Paul A.G. Sivillotti, PDCS 2005

Observation: Overconstraint

- Consistency only requires RLS before first *after* receive
 - Could be "lazy" and delay RLS
- Potential advantages:
 - Fewer messages in transit
 - less storage required to record channel state
 - Some flexibility in choosing when to RLS
 - choose less critical time
 - choose time when local state is smaller

Paul A.G. Sivillotti, PDCS 2005

Taxonomy of Laziness

- Recording Local State
 - When first **marker** received
 - Before first send or receive (after rec'd marker)
 - Before first send or *after* receive (after rec'd marker)
 - Before first *after* receive
- Sending out markers
 - When the first **marker** received
 - When local state is recorded
 - Before the first *after* send

Paul A.G. Sivillotti, PDCS 2005

Orthogonal Dimensions?

Laziness in recording local state

	I	II	III	IV
A				
B				
C				

Laziness in sending markers

Paul A.G. Sivillotti, PDCS 2005

Example 1: Chandy-Lampert

Laziness in recording local state

	I	II	III	IV
A				
B				
C				

Laziness in sending markers

Paul A.G. Sivillotti, PDCS 2005

Example 2: Prescience (A-IV)

Laziness in recording local state

	I	II	III	IV
A				
B				
C				

Laziness in sending markers

Paul A.G. Sivillotti, PDCS 2005

Example 2: Prescience (A-IV)

Paul A.G. Sivillotti, PDCS 2005

Taxonomy of Laziness

Laziness in recording local state

	I	II	III	IV
A				
B				
C				

Laziness in sending markers

Paul A.G. Sivillotti, PDCS 2005

Analysis: Flexibility

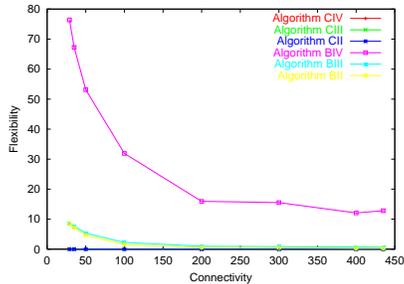
Laziness in recording local state

	I	II	III	IV
A				
B				
C				

Laziness in sending markers

Paul A.G. Sivillotti, PDCS 2005

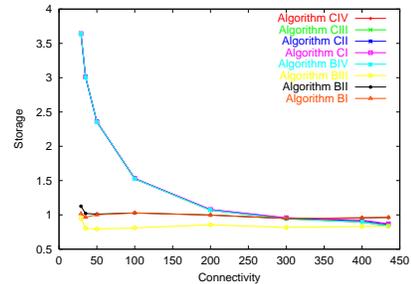
Results: Flexibility



Paul A.G. Sivilotti, PDCS 2005

19

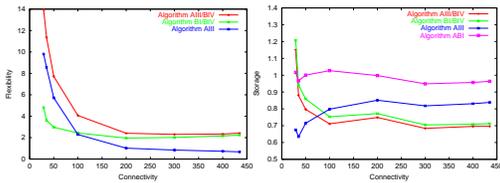
Results: Storage



Paul A.G. Sivilotti, PDCS 2005

20

Results: Hybrid Algorithms



Paul A.G. Sivilotti, PDCS 2005

21

Conclusions

- Taxonomy with 2 dimensions
 - Laziness in recording local state
 - Laziness in propagating markers
- Dimensions are not independent
 - Chandy-Lamport: AB-I
 - Unimplementable (without prescience): A-IV
- Performance evaluation
 - Flexibility and storage complexity
- Hybrid: blending A-III / B-IV
 - Dynamic mixing based on prediction of future communication events

Paul A.G. Sivilotti, PDCS 2005

22

Lazy Snapshots: A Taxonomy and Performance Study

Zhijun Liu and Paul Sivilotti

Computer Science and Engineering
The Ohio State University

paolo@cse.ohio-state.edu