

Distributed Snapshots: Determining Global States of Distributed Systems

K. Mani Chandy and Leslie Lamport

Overview

This paper by K. Mani Chandy and Leslie Lamport proposes an algorithm for determining the global state of a distributed system. The paper begins by defining the problem of creating a consistent snapshot of a distributed system. It uses the example of multiple photographers capturing a panorama in such a way that the photographs they all take, when put together create a consistent image. The problem of stable property detection is discussed as a problem whose solution lies in consistent global snapshots.

The paper then devotes a section to clearly elucidating definitions pertaining to its model of a distributed system and the associated components including processes, channels, messages, sequences, states, events and global states. The paper then presents two examples to demonstrate this model of a distributed system in action, one where everything works as intended and the other shows the problems that non-deterministic computations pose for a distributed system.

The paper then presents its algorithm for recording consistent snapshots of a distributed system, a marker-based algorithm that determines the optimal time for a process to record its state, so that it is consistent with the states recorded by other processes. The algorithm produces S^* , a recorded state which may or may not be present in the actual execution of the processes, but the paper proves that S^* is reachable from the states of execution in the program and that the termination state of the program is reachable from S^* . Finally, the paper demonstrates how stability detection can be performed with the help of consistent global states.

Discussion of Concepts Presented by the Paper

The paper clearly and effectively states the problems it wishes to solve, i.e. the problem of capturing a consistent global snapshot in a distributed system and then devises a clear and concise solution to this problem in the form of markers. This algorithm is an effective workaround to the fact that all the processes cannot take snapshots at the same time. The stitched image of each process' snapshot is a near approximation. The paper is also bookended by the stability detection problem, in the beginning by the definition of the problem and its implications, and at the end, how the problem can be solved with global distributed snapshots.

Another significantly positive aspect of the paper is the definition of the models of the distributed system. By clearly defining each aspect of the distributed system in unambiguous terms, the authors have defined a highly modular model for a distributed system. It lays the groundwork for future efforts to improve or modify the system, since the layout of the model makes it easy to substitute components.

The paper effectively demonstrates the utility of the global snapshot by proving that it is reachable from the starting state and that the termination state can be reached by this recorded state. This shows that even if the global state is not present in the actual sequence of states during the execution, it has utility.

However, this algorithm is not without its limitations. In my opinion, the model on which the algorithm is built on is based very big assumptions i.e. that channels have infinite buffers, the channels are error free and messages are received in the order they are sent. The latter two seem to be very strong assumptions to make, given that asynchrony in a distributed system is one of its most defining characteristics. In real-life distributed systems, there is often no guarantee that channels are reliable or that messages are received in the order that they are sent, throwing this algorithm's applicability into doubt.

A similar drawback is that, personally, I have to doubt the applicability of the recorded state S^* in real-life distributed systems. To tie this into the panorama photograph example in the paper, if the stitched-together photograph is of a scene that didn't actually take place, but could have possibly happened, beyond ensuring that nothing has gone horribly wrong in recording the snapshot, (e.g. the panorama photograph mistakes brought up in class), and its application in stability detection, other applications don't seem evident to me.

Despite dealing with distributed systems in which failure detection and fault tolerance are extremely important aspects for ensuring the correct functioning of the system, this paper does not contain any explanation of how to reconcile failures or measures to be taken if the event is in a state outside the invariant states. For such a tightly-defined and focused paper, in my opinion, it is peculiar that the authors do not even mention that what actions need to be taken vis-à-vis snapshots in the event of a failure lies beyond the scope of the paper.

Observations and Thoughts about Other Aspects of the Paper

As I mentioned earlier, I appreciate the clear definitions of the concepts in the paper. There are no ambiguities in terms of what each aspect of a distributed system means both as a self-contained aspect and its role in the larger system. I also appreciate the use of examples to illustrate the concepts being presented in the paper. In my opinion, although the authors could have simply stated their definitions before moving on to explaining the algorithm, the time and space spent on the examples was a good choice, as it helped me understand the concepts the authors were conveying in a succinct way.

Another aspect of this paper that I would like to point out is the minimal number of citations in this paper. In my opinion, this is an encapsulation of the paper's strengths and weaknesses. It underscores the paper's originality and self-contained nature, but also highlights the sparseness of details when it comes to applications of the algorithm in real-world distributed systems.

In terms of the presentation of the content, although the examples present in the paper illustrated the concepts extremely well, I wish the authors had included more examples or made a mention of possible edge cases, because as it stands, the paper discusses the optimal functioning of the algorithm. Non-determinism is accounted for by definition, but illustrations of how the algorithm

successfully handles rare or niche situations would have been more informative. Jumping off this point, I also think it would have been good to have an illustration of how the marker algorithm works in a larger distributed system, with more processes, or channels or tokens. Another point of discussion potentially is this: what other problems like the stability detection problem can benefit from a consistent global snapshot?

Conclusion

This paper presents an algorithm for obtaining a snapshot of the global state of a distributed system and explains its utility and application to the stability detection problem. Although the paper has certain limitations, the biggest being its assumptions of error-free channels, its clear definitions and modular nature lays the foundation for future works to overcome these drawbacks.