

The Formal Verification of Vector Clocks

Jeffrey Cheng

Advised By: Andrew Appel

April 28, 2024

Abstract

In distributed systems, events occur in parallel across numerous processes, and each process in the network must agree on the order in which these events occur. Vector clocks address this problem by assigning sortable timestamps to each event to establish a universal ordering that preserves the causal relationship between events. This preservation of causality is the key invariant that vector clocks aim to maintain. This project uses Coq to provide a formal, machine-checked proof that shows the vector clock algorithm upholds this invariant. It also demonstrates how the validity of this invariant provides correctness guarantees in a real-world system that employs vector clocks, namely distributed database operations.

1 Introduction

A distributed system is a network of computers that process events locally and communicate with each other to achieve a computational goal. However, with so many events occurring in parallel, different machines will order the events differently due to randomized latencies in communication. This can have detrimental consequences, such as copies of a database executing reads and writes in different orders.

One proposed solution to this situation is vector clocks. Vector clocks aim to order events in a manner that preserves causality, which is the idea that some event A “could have caused or affected” another event B. For example, if a process writes the value 30 and another process reads that 30, the write event “could have effected” the read event, and the read should be ordered after the write [1].

One of the original papers to propose the vector clock algorithm was a paper by Colin Fidge in 1988 [1]. In it, he provides a semi-formal, partial proof of the algorithm. To guarantee the validity of vector clocks, the goal of this project was to formally verify the vector clock algorithm in Coq, an automated proof assistant and checker. Coq only accepts proofs that abide by a set of universally accepted axioms, meaning all Coq-validated proofs, such as the one in this paper, are completely correct.

2 Background

2.1 Vector Clocks

Distributed Systems Vector clocks are used in distributed systems, where events occur on parallel processes, as depicted in Figure 1.

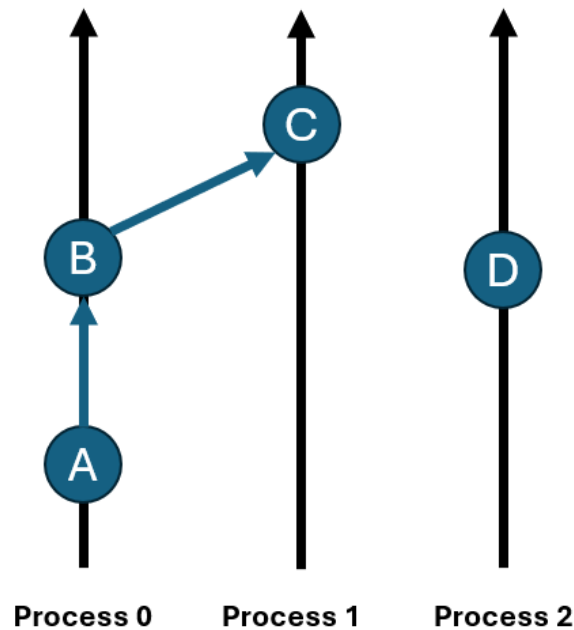


Figure 1: A example diagram of lettered events occuring on numbered processes.

In Figure 1, events A, B, C, and D happen accross processes 0, 1, and 2. These processes run independently from each other, but can transmit messages to communicate. In these scenarios, we consider three types of events:

1. Singleton Events (Event A): Events that occur locally on a process with no other effect
2. Send Message Events (Event B): Events that send a message to one or several processes
3. Receive Message Events (Event C): Events that receive a message from a different process

These three event types generalize quite well into most distributed system computations, and are the focus of vector clocks.

Causality Within these distributed computations, there is a notion of “causality.” For example, in Figure 1, it would make sense to say that event A “could have caused” or “could have affected” event B. This is because event A occurs before event B on the same process, so anything done at

A is visible to B. Similarly, it would make sense to say that event B “could have caused” event C, because B sends a message that is directly received by C. Lastly, it would make sense to say that event A “could have caused” event C, because if event A could have caused event B and event B could have caused event C, then event A could have caused event C. These three relations are the three types of causality addressed by vector clocks, which are usually denoted with an arrow:

1. Sequential Causality ($A \rightarrow B$): A occurs before B on the same process
2. Transmission Causality ($B \rightarrow C$): B sends a message that C receives
3. Transitive Causality ($A \rightarrow C$): $A \rightarrow B$ and $B \rightarrow C$

On the other hand, it would not make sense to say that event C “could have caused” event D because these events occur on two different processes that do not communicate with each other in this situation. Thus, D would have no knowledge of C.

Vector Clocks The goal of vector clocks is to assign timestamps to each event in a way that maintains causality. Specifically, the vector clock algorithm aims to uphold the following invariant:

$$A \rightarrow B \iff \text{timestamp}(A) < \text{timestamp}(B)$$

This invariant means that for Figure 1, $\text{timestamp}(A) < \text{timestamp}(B)$, $\text{timestamp}(B) < \text{timestamp}(C)$, and $\text{timestamp}(A) < \text{timestamp}(C)$. Conversely, this invariant implies that $\text{timestamp}(C) \not< \text{timestamp}(D)$ and $\text{timestamp}(D) \not< \text{timestamp}(C)$.

Vector Clock Algorithm This paper uses the vector clock algorithm proposed by Colin Fidge in 1988 [1]. This section attempts to provide a brief, intuitive understanding, but consult the original paper for details. Note that the terms “timestamp,” “vector time,” and “clock time” are used interchangeably.

In accordance with its name, timestamps in the vector clock algorithm are represented by vectors of length n , where n is the number of processes in the distributed system of interest, and each index i of any timestamp is associated with process number i . Each process keeps track of its own vector timestamp, which can be thought of as that process's current clock time that gets updated as events occur. All clock times begin as all zeroes, and when any event occurs on a process p , p updates index p of its clock by 1. When a message transmission occurs from process p to process q , the send event of the transmission attaches its timestamp to the message, and process q updates each index i of its timestamp to $\max(q_clock[i], msg_timestamp[i])$. Consider the example below.

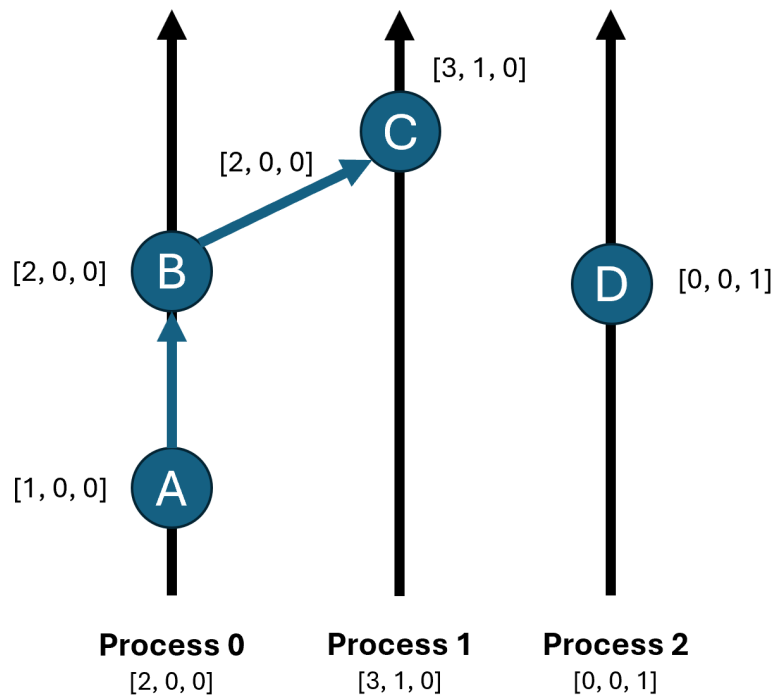


Figure 2: The distributed computation from Figure 1 with timestamps.

When event A occurs on process 0, $clock_time(process\ 0)[index\ 0]$ increases by 1, and again for event B. Event B attaches its timestamp to the message, and process 1 updates each index of its clock to the message timestamp's respective value if it is larger, giving $[2, 0, 0]$. It then increments the indices of the send and receive processes to $[3, 1, 1]$. Process 2 has no idea that this is occurring

because it has not received any messages, so it only has event D with a clock value of $[0, 0, 1]$. The vector times at the bottom represent the clock times of each process after this distributed computation has finished.

To compare the vector timestamps two events, the vector clock algorithm asserts the following invariant:

$$A \rightarrow B \iff \text{timestamp}(A)[p] < \text{timestamp}(B)[p] \quad \text{where } A \text{ occurs on process } p$$

For example, $A \rightarrow B$ and $\text{timestamp}(A)[0] < \text{timestamp}(B)[0]$, $A \rightarrow C$ and $\text{timestamp}(A)[0] < \text{timestamp}(C)[1]$, and $D \not\rightarrow C$ and $\text{timestamp}(D)[2] \not< \text{timestamp}(C)[2]$. From this point on, this proposition is referred to as the vector clock invariant, and is the main subject of interest for this paper.

Vector Clock Correctness Intuition Figure 3 below demonstrates the intuition behind why the vector clock algorithm works.

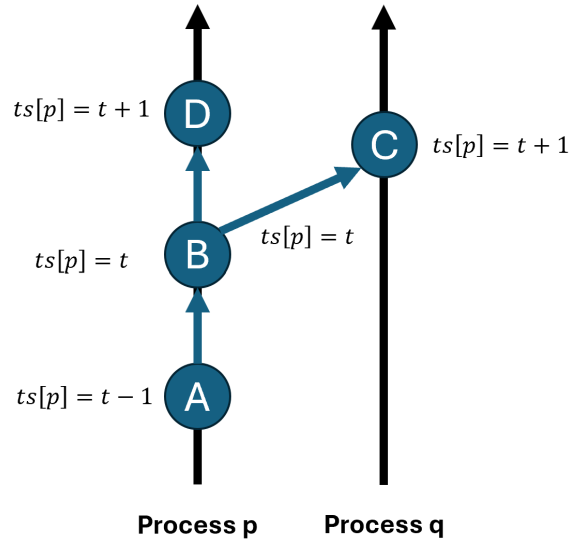


Figure 3: A distributed computation between processes p and q with timestamps abbreviated as ts .

Events A, B and D occur on process p , and their timestamps increase accordingly at index p . Event B sends a message to event C on process q , so event C “learns about” events A and B, but not D. The receive event timestamp calculation, referring to the *max* operations, is arranged so that event C’s timestamp at index p is greater than that of every event from B and before, but no event afterward. A vectorized timestamp allows each process to track this relationship with each other process. The events that C learns about are the events that “could have caused” C, which is intuitively why the vector clock invariant is true. However, we would not like to believe in vector clocks based on intuition alone. Thus, this project’s goal was to provide a formal, machine-checked Coq proof that shows that the vector clock algorithm maintains the vector clock invariant.

Coq Coq is an automated proof checker and assistant for proving theorems about functional programs. In other words, Coq allows you to write functional programs, state theorems about those programs, and apply a series of proof steps to show that your theorems are correct. As you step through your proof, Coq will tell you exactly which steps are correct and incorrect using a set of universally accepted axioms, pre-proven Coq lemmas, and any lemma that is stated and proven by the user. See Figure 4 below for an example.

```

(* Algorithm Code *)
Definition step (s : state) (e : event) : state :=
  let i := fst (fst s) in
  let c := snd (fst s) in
  let event_times := snd s in
  match e with
  | (Event p) | (Event_send p) =>
    let clocks' := update_clock c p in
    let event_times' := update_event_time p i event_times clocks' in
    ((i+1), clocks', event_times')
  | (Event_receive p s_p s_n) =>
    let timestamp_receive := event_times s_n in
    let clocks' := update_clock_receive p s_p timestamp_receive c in
    let event_times' := update_event_time p i event_times clocks' in
    ((i+1), clocks', event_times')
  end.
Definition run_computation (comp : computation) (s : state)
: (nat * clocks * time stamps) % type := fold left step comp s.

(* Lemmas *)
Lemma comp_step:
  forall (comp : computation) (tl : event),
  let event_timestamps_tl1 := run_computation (comp ++ tl :: nil) init_state
  let event_timestamps_tl2 := step (run_computation comp init_state)
  event_timestamps_tl1 = event_timestamps_tl2.
Proof.
  intros comp tl. simpl.
  apply (fold_left_app step comp (tl :: nil) init_state).
Qed.

```

```

1 goal
comp : computation
tl : event
-----
run_computation (comp ++ tl :: nil) init_state =
step (run_computation comp init_state) tl

```

Figure 4: A screenshot of the CoqIDE.

The top left quadrant of Figure 4 is part of the vector clock algorithm, implemented in Coq’s

functional programming language. In the bottom left quadrant is a lemma named “comp_step.” It asserts:

Lemma (comp_step). *For any computation $comp$ and event tl , we have:*

$$run_computation(comp ++ tl :: nil, init_state) = step(run_computation(comp, init_state), tl)$$

In other words, the two function calls, despite appearing different, always return the same value. In the top right quadrant, Coq will list defined variables, hypotheses, and your current manipulation of the theorem to prove.

In Fidge’s original 1988 paper, he provides a semi-formal proof of

$$A \rightarrow B \Rightarrow \text{timestamp}(A)[p] < \text{timestamp}(B)[p]$$

where semi-formal means an argument in English and mathematical notation. This project provides a formal proof that is verified to be correct by Coq of

$$A \rightarrow B \iff \text{timestamp}(A)[p] < \text{timestamp}(B)[p]$$

which is the full vector clock invariant.

3 Approach and Implementation

To complete this proof, there were five milestones to achieve in Coq:

1. Defining vector clocks
2. Implementing the vector clock algorithm
3. Stating the vector clock invariant

4. Proving the vector clock invariant
5. Proving a real-world vector clock guarantee

The code for all five of these milestones can be found in [this GitHub repository](#).

3.1 Defining Vector Clocks

The first step of the proof was to define all the elements of distributed systems and vector clocks in Coq's functional programming language.

3.1.1 Distributed Systems Elements

Processes and Events The code below shows how processes and events are represented in Coq.

Definition *process* := nat.

Inductive *event* : Type :=

| *Event* (*proc* : *process*)
 | *Event_send* (*proc* : *process*)
 | *Event_receive* (*proc* : *process*)
 (*send_event_proc* : *process*) (*send_event_num* : *event_index*).

The *process* datatype was defined to be a natural number representing a process ID number. The event datatype was defined to have three constructors representing each type of event:

1. Singleton Event: A local event that occurs on a process *proc*
2. Send Event: An event that sends a message from a process *proc*.
3. Receive Event: An event on process *proc* that receives a message sent by an send event *send_event_num* on the process *send_event_proc*

In Fidge’s original paper, there is a one-to-one mapping between send and receive events [1]. The definition above only requires that every receive event has exactly one send event, making it more general than the original. The proof supports the broader definition without extra work, which allows for multicasts within the distributed system, meaning one event can send the same message to multiple other processes.

There is also a slight redundancy in this definition. Each receive event stores its respective send event and send process. The send event itself stores its own process, so having the receive event store the send process as well is slightly redundant. This decision was made to simplify the proof process and requires a well-formedness guarantee that is defined in section 3.3.

Computations Vector clocks work on a distributed system with a set of events occurring over a set of processes. These situations are defined as a list of events, called a *computation*.

Definition computation := list event.

The event constructors do not provide a sufficient method of uniquely identifying events. For example, if there were two singleton events that occurred on the same process, their constructors would be exactly the same in a computation. As a result, events are identified, stored, and compared using an index into a computation, defined using the datatype *event_index*.

Definition event_index := nat.

A computation is interpreted to be a set of events that occur on a distributed system in the order that they appear in the list. For example, [Event 1, Event 3, Event_send 1, Event_receive 2 1 0] would be the distributed computation in Figure 1. Note that since event D (Event 3) has no relationship to the other events, moving Event 3 around the list would represent the same distributed computation.

3.1.2 Causality

Having defined events in Coq, we must now define what it means for $A \rightarrow B$ given two events A and B . Recall the three types of causality: sequential causality, transmission causality, and transitive causality.

Sequential Causality Given a computation, an event A , and an event B , we define B to come sequentially after A if the events occur on the same process and B comes later than A in the given computation, embodying sequential causality. To define this in Coq, we define a function `seq` that accepts a computation `comp`, an index A , and an index B that index to two events A and B in `comp`. The function returns a proposition stating that events A and B occur on the same process and A happens before B in the `comp`. Propositions in Coq are simply statements that can be true or false, and the proposition given by `seq comp A B` is true when B comes sequentially after A and false when B does not come sequentially after A .

Definition `seq (comp : computation) (A B : event_index) : Prop :=`
`(event_process comp A = event_process comp B) ∧ (A < B).`

Transmission Causality Given a computation, an event A , and an event B , we define B to be the receive event of a send event A if B is a receive event, A is the send event stored in B 's constructor, and the send process number in B 's constructor is A 's process number, embodying transmission causality. We define a Coq function `msg` that again accepts a computation `comp`, an index A , and an index B that index to events A and B in `comp`. It returns a proposition that is true when $A \rightarrow B$ by transmission causality and false otherwise.

Definition `msg (comp : computation) (A B : event_index) : Prop :=`
`match nth B comp (Event 0) with`
`| Event_receive _ s_p s_n ⇒ A = s_n ∧ (event_process comp s_n = s_p)`
`| _ ⇒ False`

end.

This function fetches the event B from `comp` using index B and pattern matches the constructor of the retrieved event. If event B is a singleton event or send event, it cannot be the case that $A \rightarrow B$ by transmission causality, so the proposition should always be false. If B is a receive event, it returns the proposition that index A matches the the index of B 's corresponding send event and event A occurs on the process that B received a message from according to its constructor. Again, the returned proposition is true when $A \rightarrow B$ by the above definition of transmission causality and false otherwise.

General Causality Lastly, given a computation, an event A , and an event B , we define the general relation $A \rightarrow B$ using three constructors, one for each type of causality.

1. $A \rightarrow B$ if you can provide evidence that B comes sequentially after A
2. $A \rightarrow B$ if you can provide evidence that B is the receive event to the event send A
3. $A \rightarrow C$ if you can provide evidence that $A \rightarrow B$ and $B \rightarrow C$ for an event B in the computation

In Coq, we refer to $A \rightarrow B$ as the “after” relation, where saying B comes “after” A is the same as saying A “could have caused” B . We define the “after” relation below.

```
Inductive after : computation → event_index → event_index → Prop :=
  | after_sp (comp : computation) (A B : event_index) (H : seq comp A B) : after comp A B
  | after_dp (comp : computation) (A B : event_index) (H : msg comp A B) : after comp A B
  | after_trans (comp : computation) (A B C : event_index) (Htrans1 : after comp A B) (Htrans2
: after comp B C) : B < length comp → after comp A C.
```

Similar to `seq` and `msg`, `after` accepts a computation `comp`, an index A , and an index B , as indicated to the right of the colon of each constructor. `after comp A B` is a proposition that is true when $A \rightarrow B$ and false otherwise. You can demonstrate that an `after` proposition is

true using one of its three constructors. First, you could provide a `comp`, `A`, `B`, and a proof that `seq comp A B` is true, to show `after comp a b` by sequential causality. Second, you could provide a `comp`, `A`, `B`, and a proof that `msg comp A B` is true, to show `after comp a b` by transmission causality. Third, you could provide a `comp`, `A`, `B`, `C`, evidence that `after comp A B` is true, evidence that `after comp B C` is true, and evidence that `B` is a valid index into `comp` to show that `after comp A C` is true.

Vector Timestamps The novelty of Fidge’s paper comes in how events are timestamped. Each event receives a timestamp that is a vector of natural numbers, indexed by process numbers. Thus, a vector timestamp is defined in Coq as a function that maps process numbers to natural numbers. Arrays were not used because a strictly functional implementation was desired, and a list was not used because a function allows for infinite processes.

Definition *time* := *nat*.

Definition *vclock_time* := *process* → *time*.

Each process maintains a vector timestamp as its clock value, so the set of all clock values is represented as a function mapping process numbers to their clock times. The vector clock algorithm aims to timestamp each event, so its output is represented as a mapping of events to their timestamp for a particular computation.

Definition *clocks* := *process* → *vclock_time*.

Definition *time_stamps* := *event_index* → *vclock_time*.

3.2 The Vector Clock Algorithm

The vector clock algorithm is a function named `run_computation` takes in a `computation` and produces a `time_stamps` function. It processes the list of events from head to tail, updating the process clocks and assigning a timestamp for each event. To make the proofs by induction

easier, this was defined as a `fold_left` with a `step` function to process each event. The code below feeds a computation `comp` to the vector clock algorithm `run_computation` and tests if $\text{timestamp}(A)[p] < \text{timestamp}(B)[p]$.

```
let p := event_process comp A in
let event_timestamps : time_stamps :=
  snd (run_computation comp init_state) in
event_timestamps A p < event_timestamps B p
```

3.3 Theorem Statement

$$A \rightarrow B \iff \text{timestamp}(A)[p] < \text{timestamp}(B)[p] \quad \text{where } A \text{ occurs on process } p$$

This is the vector clock invariant that we want to prove. Below is the statement of the invariant in Coq.

Theorem *order* :

```
∀ (comp : computation) (A B : event_index) (p : process),
A < length comp →
B < length comp →
p = event_process comp a →
(∀ (i : event_index), well_formed comp i) →
after comp A B ↔
let event_timestamps : time_stamps :=
  snd (run_computation comp init_state) in
event_timestamps A p < event_timestamps B p.
```

There are four assumptions required for the invariant. Events A and B must be in the given computation, p must be the process that A occurs on, and the computation must be well-formed. There are two well-formedness rules that are not encapsulated by the definitions of the vector

clock elements and must be assumed as hypotheses. First, every receive event must come after its corresponding send event. Second, the send process number stored in a receive event constructor must be equal to the process stored by the corresponding send event's constructor. Each of these well-formedness statements are written as Coq propositions, and the `well_formed` proposition that asserts that both the well-formedness propositions are true for some event in `comp`.

3.4 Theorem Proof

To prove the theorem, both the forward and backward direction of the if-and-only-if had to be proven. The forward direction was broken into three subcases, one for each type of `after` constructor to show that the invariant holds for all three types of causality. From there, each subcase was proven using induction on the length of a list. This meant that the proof goal could be assumed as the induction hypothesis for any list of events of length n , and it had to be shown that appending an event to the end of the list would not break the invariant. Since the event timestamps assigned to events do not change when an event is appended to the end of the list, our induction hypothesis for lists of length n is quite powerful when proving the invariant for lists of length $n + 1$. If A nor B are the tail element of the list, we can remove the tail element, and our induction hypothesis proves the invariant for A and B. This meant that each one of the subcase proofs only had to deal with the scenario where A or B is the tail element of the list.

3.4.1 The Forward Direction: Sequential Events

$$seq\ comp\ A\ B \rightarrow event_timestamps\ A\ p < event_timestamps\ B\ p$$

The first subcase to prove is that if B occurs sequentially after A on the same process, then $timestamp(A)[p] < timestamp(B)[p]$, where A occurs on process p . For each event that occurs on a process p , the vector clock algorithm increases p 's clock time at index p by 1. Since A and B occur on the same process and B occurs after A, it is quite straightforward to see how index p of

B's timestamp would be larger than A's.

3.4.2 The Forward Direction: Message Transmission Events

$$\text{msg_comp } A \rightarrow B \rightarrow \text{event_timestamps } A \ p < \text{event_timestamps } B \ p$$

The second subcase to prove is that if B receives a message that A sends, $\text{timestamp}(A)[p] < \text{timestamp}(B)[p]$, where A occurs on process p . According to the vector clock algorithm, $\text{timestamp}(B)[p] = \max(\text{local_clock}[p], \text{timestamp}(A)[p] + 1)$. No matter which value is larger, $\text{timestamp}(A)[p]$ is assigned a value larger than $\text{timestamp}(A)[p]$.

3.5 The Forward Direction: Transitive Events

$$\text{after_comp } A \rightarrow B \rightarrow \text{after_comp } B \rightarrow C \rightarrow \text{event_timestamps } A \ p < \text{event_timestamps } C \ p$$

The last subcase of the forward direction is that if B comes after A and C comes after B, $\text{timestamp}(A)[p] < \text{timestamp}(C)[p]$ where A occurs on process p . Given the first two subcases, this seems quite straightforward, but there is a catch. When checking if $A \rightarrow B$, the index p used to index the timestamps is always the process that the A, the first event, occurs on. For example, according to the invariant, $A \rightarrow B \Leftrightarrow \text{timestamp}(A)[p] < \text{timestamp}(A)[p]$ and $B \rightarrow A \Leftrightarrow \text{timestamp}(B)[q] < \text{timestamp}(A)[q]$ where A occurs on process p and B occurs on process q . By comparing at different indices based on the direction of the arrow, we can have $A \rightarrow B$ and $B \rightarrow A$ without being forced to assign timestamps of equal values.

What this means for the transitive property is that our invariant gives us

$B \rightarrow C \Rightarrow \text{timestamp}(B)[q] < \text{timestamp}(C)[q]$, where B occurs on process q and q might not be equal to p . As a result, $B \rightarrow C$ does not guarantee $\text{timestamp}(B)[p] < \text{timestamp}(C)[p]$, which we wanted for the transitive inequality. First, it must be proven that

$A \rightarrow B \Rightarrow \text{timestamp}(A)[p] \leq \text{timestamp}(B)[p]$ for all i , which is proven in essentially the same

manner as the forward direction of the original invariant. This lemma does have a trivial transitive subcase.

Once the less than or equal to lemma is proven, we can reason that $A \rightarrow B \Rightarrow \text{timestamp}(A)[p] < \text{timestamp}(B)[p]$ and $B \rightarrow C \Rightarrow \text{timestamp}(B)[p] \leq \text{timestamp}(C)[p]$, yielding $\text{timestamp}(A)[p] < \text{timestamp}(C)[p]$ where A occurs on process p .

3.5.1 The Backward Direction

The backward direction of the vector clock invariant is the more interesting result. To prove it, we prove the contrapositive, meaning we show that

$$A \nrightarrow B \Rightarrow \text{timestamp}(A)[p] \geq \text{timestamp}(B)[p]$$

which is logically equivalent to $\text{timestamp}(A)[p] < \text{timestamp}(B)[p] \Rightarrow A \rightarrow B$. Again, we use induction on the length of a list, allowing us to assume $A \nrightarrow B \Rightarrow \text{timestamp}(A)[p] \geq \text{timestamp}(B)[p]$ for any events A, B in a list of length n . We have to show that this invariant holds when an element is appended to the end of the list.

Like before, if neither A nor B are the tail element, the inductive hypothesis makes this case trivial. The tail element on a process p is always the last event for the index p . Thus, if A is the tail element on a process p , then $\text{timestamp}(A)[p]$ cannot be smaller than any other event's timestamp at index p . As a result, $\text{timestamp}(A)[p] \geq \text{timestamp}(B)[p]$ for all other events B. The intuition for B as the tail element is more complicated.

Case 1 First take the case where B is a receive event that inherited index p of its timestamp from a previous send event X during the $\max(\text{local_clock}[p], \text{msg_timestamp}[p])$ calculation. Since $X \rightarrow B$ and $A \nrightarrow B$, it must be that $A \nrightarrow X$, otherwise $A \rightarrow B$ by transitive causality. X precedes B, so it is not the tail element and we can use it with the inductive hypothesis: $A \nrightarrow X \Rightarrow$

$\text{timestamp}(A)[p] \geq \text{timestamp}(X)[p]$. To calculate the relation between $\text{timestamp}(X)[p]$ and $\text{timestamp}(B)[p]$, we do case analysis on all the placements of X as shown below:

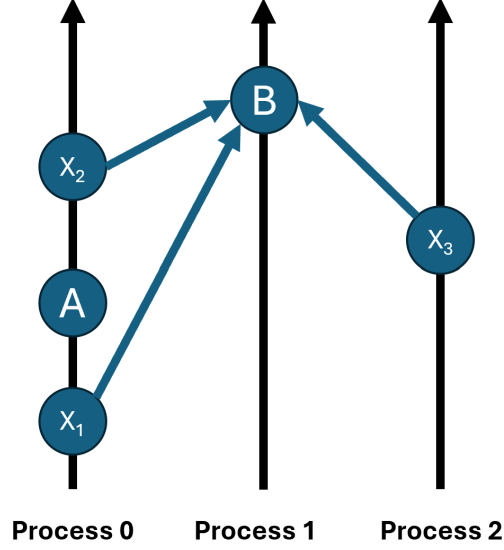


Figure 5: All placements of a send event X for a receive event B .

For X_1 , $\text{timestamp}(X_1)[p] + 1 = \text{timestamp}(B)[p]$ and $\text{timestamp}(X_1)[p] < \text{timestamp}(A)[p]$, yielding $\text{timestamp}(A)[p] \geq \text{timestamp}(B)[p]$. For X_2 , X_2 comes sequentially after A , meaning $A \rightarrow X$, which is a contradiction. For X_3 , $\text{timestamp}(X_3)[p] = \text{timestamp}(B)[p]$, giving $\text{timestamp}(A)[p] \geq \text{timestamp}(B)[p]$. These three cases are exhaustive because X_1 and X_2 consider the cases where A and X are on the same process, and X_3 generalizes to any case where they are not.

Case 2 Now take any other case. This means that B is either a singleton event, a send event, or a receive event who inherited index p of its timestamp from its own process's clock during the $\max(\text{local_clock}[p], \text{msg_timestamp}[p])$ calculations. Let X be the last event on the same process as B excluding B . This means that B happens immediately after X on the same process, so $X \rightarrow B$. We now employ a similar argument as case 1.

Since $X \rightarrow B$ and $A \nrightarrow B$, $A \nrightarrow X$. By the inductive hypotheses, $A \nrightarrow X \Rightarrow \text{timestamp}(A)[p] \geq \text{timestamp}(X)[p]$. Since B comes immediately after X on some process q , q is the only timestamp index that would increase from X to B, so $\text{timestamp}(B)[p] = \text{timestamp}(X)[p]$. Note that q cannot be equal to p , otherwise B would come sequentially after A and $A \rightarrow B$, which is a contradiction. We can now say that $\text{timestamp}(A)[p] \geq \text{timestamp}(B)[p]$. This concludes the final subcase for proving the main vector clock theorem.

3.6 Real-World Application

Having proved the main vector clock invariant, we can show how useful this guarantee is with a real world algorithm that employs vector clocks. One such application of vector clocks is distributed database operations. For this case, we can think of each process being an independent user and each event being a database operation that a user executes. Certain database operations follow directly from others. However, other sets of operations can be executed independently of each other because their users were not in communication. We can use vector clocks to know exactly which operations can affect other operations.

For simplicity, we define our database's state to be a single natural number. We consider every event, including send and receive events, to be a database increment and read, incrementing the numerical state and assigning that state value to the event. Singleton and send events simply increment the state by one, but receive events sum the states of the send process and receive process and then increment by 1 to represent a merging of states. That way, changing events before either the receive event or its corresponding send event will modify the merged state value. This is akin to executing a database operation on a database, then viewing the new state that you created at that point in time. As a result, modifying the list of database operations that occur may or may not modify certain operations' assigned state.

We want to reason about how modifying one operation affects others. One way of modifying an operation is stifling the state update of an operation A, then seeing whether that affects the

state value of an operation B. The state calculation function takes in one index and prevents the operation at that index from incrementing the state. If that operation is a send or receive event, the message transmission still works, just without a state increment or summation upon reception.

In the real world, this function would likely use vector clocks to tell when stifling the state update of operation A could affect an operation B. The developer would reason that $A \rightarrow B$ is synonymous with “stifling A affects B”, and our vector clock invariant would allow them to say $\text{timestamp}(A)[p] < \text{timestamp}(B)[p]$ if and only if stifling A affects B. This allows them to use vector clocks to tell when certain operations affect others. Thus, our distributed database invariant is as follows:

Theorem *affects*:

$$\begin{aligned} & \forall (comp : computation) (A B : event_index) (p : process), \\ & A < length\ comp \rightarrow \\ & B < length\ comp \rightarrow \\ & p = event_process\ comp\ A \rightarrow \\ & (\forall (i : event_index), well_formed\ comp\ i) \rightarrow \\ & after\ comp\ A\ B \leftrightarrow \\ & let\ op_states := apply_operations\ comp\ state\ (length\ comp)\ in \\ & let\ op_states_rem_A := apply_operations\ comp\ state\ A\ in \\ & op_states\ b \neq op_states_rem_A\ b. \end{aligned}$$

For this proof and most other real world application proofs, the intuitive approach would be to show that $A \rightarrow B \iff op_states\ b \neq op_states_rem_A\ b$. However, since the state calculation algorithm mirrors the vector clock algorithm quite closely, the proof was done by showing $\text{timestamp}(A)[p] < \text{timestamp}(B)[p] \iff op_states\ b \neq op_states_rem_A\ b$. Both these statements are equivalent, and the final result is

$$A \rightarrow B \iff op_states\ b \neq op_states_rem_A\ b \iff \text{timestamp}(A)[p] < \text{timestamp}(B)[p]$$

The proof of this theorem used similar techniques to the proof of the main vector clock theorem. Using induction on the length of the list, we only really need to worry about the case where B is the tail of the list. Find the event X that B inherits index p of its timestamp from. $\text{op_states } x \neq \text{op_states_rem_A } x \iff \text{timestamp}(A)[p] < \text{timestamp}(X)[p]$ by the induction hypothesis, and manually calculate the relationship between X and B for both their timestamps and database states, given B builds exactly one step off of X.

4 Results and Conclusion

From this project, we are now able to conclude the validity of vector clocks with 100% certainty. By providing Coq with incredibly detailed proof steps, we have shown that there is no possible corner case in the algorithm that was forgotten nor incorrectly reasoned about.

This provides a pleasant guarantee of correctness for those that employ vector clocks in their systems, such as distributed databases. By assigning vector timestamps to the events of a distributed system, developers can know exactly when one event can affect another. This allows them to remove, add, or change events to distributed computations and know exactly which states will be affected. For example, database users could delete database operations and know which database reads will come back differently to avoid unwanted consequences. GitHub users could delete commits and know which future commits will be affected. We have also formally proven that these systems express causality exactly as defined. In other words, $A \rightarrow B$ truly is synonymous with “stifling A’s state update affects B’s state”, which is an intuitive but non-trivial idea. This shows that vector clocks are the correct tool of choice for these applications.

For most of the computer science world, semi-formal proofs are often enough to convince people of correctness, and vector clocks are commonly used without having a formal proof until now. However, this project does not just provide a proof to a widely accepted notion. There are many vector clocks optimizations whose correctness are not as obvious. To show that these

optimizations are not too greedy and do not violate the vector clock invariant, this formal proof is a good foundation and model of how to show correctness in algorithms related to vector clocks. One straightforward corollary is the correctness of Lamport clocks, which only address the forward direction of the vector clock algorithm. Ultimately, this project provides a foundation and reference for many proofs to come, while allowing those who employ vector clocks in their own systems to sleep easier at night.

5 Acknowledgements

I would like to thank my advisor, Professor Andrew Appel, for his guidance and feedback throughout this project.

I pledge my honor that this paper represents my own work in accordance with University regulations. - Jeffrey Cheng

References

- [1] Colin Fidge. Timestamps in message-passing systems that preserve the partial ordering. *Australian Computer Science Communications*, 10(1):56–66, 1988.