

# SNOW Revisited and Bounded Latency READ Transactions

## Abstract

Large-scale web services are built on distributed storage systems, that *shard* (i.e., partition) data across many machines within a datacenter and *geo-replicate* the data across several datacenters. At the heart of designing such systems is the trade-off between the strength of guarantees provided to overlying applications and the performance of the system. The trade-off due to sharding, in contrast, is far less understood with only one result that is relevant for the READ transactions that dominate real-world workloads. This result is the SNOW Theorem [14] that proves the ideal READ transactions that have optimal latency and the strongest guarantees—i.e., “SNOW” READ transactions—are impossible. Our work builds on the SNOW Theorem to clarify this fundamental trade-off by closing several open questions and provide rigorously proved results. These questions relate to the possibility of SNOW with two clients, multiple writers and a single reader, and when client-to-client communication is allowed. We prove that SNOW is possible in some of these settings, but that it is impossible in the most general case even with client-to-client communication. We present the first algorithms with bounded latency for READ transactions with the strongest guarantees. We provide rigorous proofs for the SNOW theorem for the three-client and show the impossibility of SNOW properties for transaction system with two-clients. Our proofs require delicate reasoning about the inter-play of the SNOW properties and asynchrony, therefore, we rely heavily on IO automata theory to derive these results.

**Keywords and phrases** impossibility results, strict-serializability, transactions, IO Automata

## 1 Introduction

Today’s web services are built on *distributed storage systems* that provide fault tolerant and scalable access to data. Distributed storage systems scale their capacity and throughput by *sharding* (i.e., partitioning) data across many machines within a datacenter, i.e., each machine stores a subset of the data. They also *geo-replicate* the data across several geographically dispersed datacenters to tolerate failures and to increase their proximity to users.

Distributed storage systems abstract away the complexities of sharding and replication from application code by providing guarantees for accesses to data. These *guarantees* include consistency and transactions. *Consistency* controls the values of data that accesses may observe and *transactions* dictate what accesses may be grouped together. Stronger guarantees provide an abstraction closer to a single-threaded environment, greatly simplifying application code. Ensuring the guarantees hold, however, often comes with worse performance. Therefore, the tradeoff between performance and correctness guarantees lies at the heart of designing such systems.

The performance-guarantee tradeoffs that result from replication have been well-studied with several well-known impossibility results [1, 7, 8, 12]. For instance, the CAP Theorem [8] proves that system designers must choose either availability during network partitions (performance) or strong consistency across replicas (guarantee). However, little prior work exists on what performance-guarantee tradeoffs result from sharding (splitting datasets across multiple servers).

Understanding the performance-guarantee tradeoff due to sharding is important because user requests are typically handled across many shards but within a single nearby datacenter (replica). This is particularly true for the reads needed to handle a user request, which



© Anonymous author(s);

licensed under Creative Commons License CC-BY

42nd Conference on Very Important Topics (CVIT 2016).

Editors: John Q. Open and Joan R. Access; Article No. 23; pp. 23:1–23:32

Leibniz International Proceedings in Informatics



LIPIC Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

are what dominate real-world workloads: Facebook reported 500 reads for every write in their TAO system [4] and Google reported three orders of magnitude more reads than general transactions for their F1 database that runs on their Spanner system [6]. In this work we focus on clarifying the performance-guarantee tradeoff for reads that results from sharding. Distributed storage systems group read requests (each to access a separate shard) into *READ transactions* that together return a consistent, cross-shard view of the system. Whether a view is consistent is determined by the consistency model a system provides. The ideal READ transactions would have the strongest guarantees: They would provide strict serializability [16], the strongest consistency model, and they could be used in a system that also includes *WRITE transactions* that group write requests (each for a separate shard) together. The alternative to the latter property are READ transactions that can only be used in systems that have non-transactional, *simple writes*.

The ideal READ transactions would also provide the best performance. In particular, they would provide the lowest possible latency because the prevalence of reads makes them dominate the user response times that are aggressively optimized by web services [5, 11, 17]. The *optimal latency* for a READ transaction is to match the latency of non-transactional, *simple reads*: complete in a single round trip of non-blocking parallel requests to the shards that return only the requested data [14].

## 1.1 Previous Results and Open Questions

To the best of our knowledge, the only result in the sharding dimension that is relevant to READ transactions is the SNOW Theorem [14]. The SNOW Theorem is an impossibility result that proves no READ transaction can provide **Strict** serializability with **Non-blocking** client-server communication that completes with **One** response per read in a system with **WRITE** transactions (§2.1). It shows there is a fundamental tradeoff between the latency and guarantees of READ transactions that system designers must grapple with, they must pick either the strongest guarantees (S and W) or optimal latency (N and O).

SNOW is trivially possible in systems with a single client or a single server because the single entity naturally serializes all transactions. The SNOW Theorem shows SNOW is impossible in systems with at least three clients and at least two servers. It explicitly leaves open the question of the possibility of SNOW in a system with two clients.

In addition, the model used in the prior work implicitly leaves open several questions. It assumed the three clients were a single writer and multiple readers (SWMR). This leaves open the possibility of SNOW with multiple writers and a single reader (MWSR). The SNOW Theorem also implicitly disallowed client-to-client communication (C2C). This leaves open the possibility of SNOW when client-to-client communication is allowed.

In this work, the new impossibility results built upon the SNOW Theorem are philosophically similar to other impossibility results—such as FLP [7] and CAP [3, 8]—in that they help system designers avoid wasting effort in trying to achieve the impossible. That is, the SNOW Theorem identifies a boundary in the design space of READ transactions, beyond which no algorithms can possibly exist. By revisiting SNOW, our work makes this boundary more precise.

## 1.2 Our Results

Our work builds on the SNOW Theorem to clarify this fundamental tradeoff by closing these open questions. First, we present a new, rigorous proof, using I/O automata [15], of the impossibility of SNOW with three or more clients even when client-to-client communication

is allowed (§4). Then, we find that SNOW's possibility with two clients depends on if client-to-client communication is allowed: when it is not allowed, SNOW is impossible (§5.1); when it is allowed, SNOW is possible (§5.2). Our algorithm demonstrating the latter possibility also shows that SNOW is possible with multiple writers and a single reader using client-to-client communication.

## 2 Transactions Processing System

Web services typically have two tiers of machines within a datacenter: a stateless frontend tier and a stateful storage tier. The frontends handle user requests by executing application logic that generates sub-requests to read/write data in the storage tier that shards (or splits) data across many machines. We refer to the front-ends as the *clients*, the storage machines as the *servers* and the stored data items as *objects*, to match common terminology. While web services are typically geo-replicated, we focus on sharding within a datacenter because the reads that dominate their workloads are handled within a single datacenter.

We consider a transaction processing system that comprises a set of read/write objects  $\mathcal{O}$ , where each object  $o \in \mathcal{O}$  is maintained by a separate server process, and also another set of processes, we refer to as *clients*, that can initiate transactions. The system allows two types of transaction: READ transaction, a group of read requests for the values stored in some subset of objects in  $\mathcal{O}$ ; and WRITE transaction, a group of write requests intending to update the values stored in some subset of objects  $\mathcal{O}$ . A read-client executes only READ transactions, while write-client executes only WRITE transactions; no client executes both types of transaction.

A typical READ transaction, we denote as  $READ(o_{i_1}, o_{i_2}, \dots, o_{i_q})$  or in short by  $R$ , consists a set of individual read requests  $read(o_{i_1})$ ,  $read(o_{i_2})$  and  $read(o_{i_q})$  to read values in objects  $o_{i_1}, o_{i_2}, \dots, o_{i_q}$ , respectively.  $read(o)$  denotes a read that intends to read the value of object  $o$ . A typical WRITE, denoted as  $WRITE((o_{i_1}, v_{i_1}), (o_{i_2}, v_{i_2}), \dots, (o_{i_p}, v_{i_p}))$  or in short as  $W$ , consists of a set of which requests to update the values of objects  $o_{i_1}, o_{i_2}, \dots, o_{i_p}$  with  $v_{i_1}, v_{i_2}, \dots, v_{i_p}$ , respectively.

A read (or write) client initiates a READ (or WRITE) transaction with an invocation step  $INV(R)$  (or  $INV(W)$ ), then it carries out the read or write operations in the transaction; and eventually completes the transaction with a  $RESP(R)$  (or  $RESP(W)$ ). After the completion of the reads or writes in a transaction the client responds, in the case of  $R$ , with the values of objects; and, in the case of  $W$ , an *ok* status, to the external client.

We assume that the network channels are reliable but asynchronous, i.e., any message sent by a process will eventually arrive at its destination uncorrupted. We assume local computations are asynchronous, i.e., local computations at various processes proceed at arbitrary and unpredictable speeds. When a client receives a transaction request, usually from an external client, such as an user's device, it executes the transaction, denote by  $R$  or  $W$ , and finally, responds to the external client with the results.

We model a distributed algorithm using the I/O Automata (see [15] for a detailed account or the appendix, for a brief description). In the rest of this paper, for any execution of an automaton  $\mathcal{A}$ ,  $\sigma_0, a_1, \dots, a_k, \sigma_k \dots$ , where  $\sigma$ 's and  $a$ 's are states and actions, we use the notation  $a_1, \dots, a_k \dots$  that shows only the actions while leaving out the states to simplify notation. In our model, an individual read, such as  $read(o)$ , in some read transaction  $R$  initiated by some read client  $r$  consists of the following sequence of actions: after  $INV(R)$  at  $r$ ,  $r$  sends a message  $m$  (requesting the value stored in  $o$ ) to a server  $s$  via action  $send(m)_{r,s}$ . When  $s$  receives  $m$  via action  $recv(m)_{r,s}$  it sends the value  $v$  (stored in  $o$ ) to  $r$  via action

137  $send(v)_{s,r}$ . Then  $read(o)$  completes as soon as  $r$  receives  $v_j$  via action  $recv(v)_{s,r}$ .

## 138 2.1 SNOW Properties

139 In this subsection, we describe the SNOW properties. The SNOW properties for a transaction  
 140 processing system can be described by requiring that any fair execution of the system satisfies  
 141 the following four properties: (i) *Strict serializability* (S), which means there is a total  
 142 ordering of the transactions such that all transactions in the resulting execution appear to  
 143 be processed by a single machine one at a time; (ii) *Non-blocking operations* (N), which  
 144 means that the servers respond immediately to the read requests of a READ transaction  
 145 without waiting for any input from other processes; (iii) *One response per read* (O), requires  
 146 that any read operation consists of one-round trip of communication with a server, and also,  
 147 the server responds with a message that contains exactly one version of the object value;  
 148 and (iv) *WRITE transactions that conflict* (W) implies the existence of concurrent WRITE  
 149 transactions that update the data objects while READ transactions are in progress reading  
 150 the same objects. Below we describe the individual properties of the SNOW properties in  
 151 more detail.

152 **Strict serializability (S).** By *strict serializability* (for a formal definition please see [10]),  
 153 we mean each WRITE or READ transaction is executed atomically, at some point in an  
 154 execution between the invocation and response events.

155 Here we describe the one-response and non-blocking property of any transaction processing  
 156 system in this work. The non-blocking and the one-response properties are essentially defined  
 157 as properties of read operations to individual object. For the purpose of elucidation, we  
 158 consider an execution  $\alpha$  of a transaction processing system  $T$ , that has a set of objects  $\mathcal{O}$ ,  
 159 where there is a READ transaction  $READ(o_{i_1}, o_{i_2}, \dots, o_{i_q})$ , in short  $R$ , invoked at some  
 160 reader  $r$ , such that  $R$  contains a read  $read(o_j)$  for some  $o_j \in \mathcal{O}$  maintained at server  $s_j$ . By  
 161  $prefix(\alpha', a)$  we denote the execution fragment in  $\alpha$  from the beginning of  $\alpha$  to the state  
 162 immediately following any action  $a$  in  $\alpha$ .

163 **Non-blocking reads (N).** The *non-blocking* property means that if  $r$  sends any message  
 164 to any  $s_i$  ( $s_i$  manages object  $o_i$ ) during the transaction then  $s_i$  can respond to  $r$  without  
 165 waiting for any external input event, such as the arrival of messages, any mutex operations,  
 166 time, etc. This property ensures that READ transactions are delayed only due to delay in  
 167 message delivery between  $r$  and  $s_i$ . We define this property formally as follows.

168 **► Definition 1 (Non-blocking read (N)).** Suppose in  $\alpha$ , following the action  $INV(R)$ , the  
 169 actions  $recv(m_j^r)_{r,s_j}$  and  $send(v_j)_{s_j,r}$  corresponding to  $read(o_j)$ , occurs at  $s_j$ . Then there  
 170 exists an execution  $\alpha'$  of  $T$  such that

- 171 (i) The execution fragments  $prefix(\alpha, recv(m_j^r)_{r,s_j})$  and  $prefix(\alpha', recv(m_j^r)_{r,s_j})$  are identical,  
 172 where  $prefix(\alpha, a)$ .
- 173 (ii) In  $\alpha'$  the action  $send(v_j)_{s_j,r}$  at  $s_j$  occurs after  $recv(m_j^r)_{r,s_j}$  without any input action in  
 174 between.

175 **One-response per read (O).** The *one-response* property requires that each read  
 176 operation,  $read(o_i)$ , during any READ transaction, completes successfully in one round of  
 177 client-to-server communication and the *one-version* states that exactly one version of the  
 178 value is sent by server  $s_i$ , that manages  $o_i$ , to  $r$ . *One-round* consists of a read request from  
 179 the client initiating the read operation to the server and the response containing value sent  
 180 by the server.

181 ► **Definition 2** (One response per read (O)). Suppose in  $\alpha$ , the action  $INV(R)$ , the actions  
 182 occur at  $s_j$ , then in  $\alpha$  there exists exactly a pair of actions  $recv(m_j^r)_{r,s_j}$  and  $send(v_j)_{s_j,r}$ ,  
 183 corresponding to  $R$ , occur at  $s_j$ , such that  $v_j$  is the object value of  $o_j$ .

184 If the reads, of some READ transaction, of a transaction processing system respect the  
 185 non-blocking and one-response properties then each read includes one-round trip from client  
 186 to server, where the server returns only the requested value as soon as it receives the request.  
 187 It is worth noting that the READ transaction can complete only after all the  $read(\cdot)$ s in it  
 188 complete.

189 **WRITE transactions that conflict (W).** The *conflicting writes* property states that  
 190 READ transactions complete even in the presence of concurrent WRITE transactions, where  
 191 the write operations might update some objects that are also being read by read operations  
 192 in READ transaction. This shows that READ transactions can be invoked at any point,  
 193 even in the presence of ongoing WRITE transactions. Note that the liveness of any WRITE  
 194 transaction is not implied by any of the SNOW properties; however, for useful practical  
 195 systems the WRITE transactions must eventually complete. Therefore, we assume that every  
 196 WRITE transaction eventually completes via the *RESP* event, and think of this constraint  
 197 as a part of the W property.

198 **The SNOW Theorem.** Consider a transaction processing system with an asynchronous  
 199 network where a set  $\mathcal{O}$  of objects are maintained by individual server processes, with at least  
 200 one write client and at least two read clients. Then the SNOW Theorem [14] can be stated  
 201 as follows.

202 “For any transaction processing system in an asynchronous setting, with at least one  
 203 writer and two reader clients, and at least two sharded objects, it is impossible to have an  
 204 algorithm such that all of its executions guarantee the SNOW properties.”

### 205 3 Technical Preliminaries

206 In this section, we describe the proof techniques and some useful preliminary results for  
 207 proving the main impossibility results. Those main results will be proven by showing  
 208 contradictions. In each proof, we use a simple system consisting of two servers,  $s_x$  and  $s_y$ ,  
 209 and either two or three clients. One of the clients is a writer  $w$ , which initiates only WRITE  
 210 transactions. One or two of the clients are readers,  $r_1$  and  $r_2$ , which initiate only READ  
 211 transactions.

212 Servers  $s_x$  and  $s_y$  store values for objects  $o_x$  and  $o_y$ , respectively. the initial values of  
 213  $o_x$  and  $o_y$  are  $x_0$  and  $y_0$ , respectively. Because there is one object on each server the server  
 214 and object identifiers are often used interchangeably to remove redundancy. For instance, we  
 215 simply say that  $s_x$  returns  $x_0$  to the client that initiated the transaction, which means that  
 216  $s_x$  returns the value  $x_0$  of object  $o_x$  at the end of the READ transaction.

217 Our proofs often use a special type of execution fragment, named non-blocking fragments,  
 218 that represent the READ transaction algorithm is non-blocking and returns one version of  
 219 each object. The one-round property is captured by allowing only one non-blocking fragment  
 220 on each server for a READ transaction. Our proof strategy plays non-blocking fragments  
 221 against the requirements of strict serializability and write isolation under the freedom of  
 222 network asynchrony. We explain non-blocking fragments and helper notations in the context  
 223 of execution  $\alpha$ , of the system described as above, as follows:

- 224 1. *Non-blocking fragments*  $F_{i,j}(\alpha)^{(v_j)}$ ,  $j \in \{x, y\}$ . For a READ transaction  $R_i$  by reader  
 225  $r_i$ ,  $i \in \{1, 2\}$ , suppose there is a execution fragment that starts with  $recv(m_j^r)_{r_i,s_j}$  and

- ends with  $send(v_j)_{s_j, r_i}$ , both of which occur at  $s_j$ . Moreover, suppose there is no other input action at  $s_j$  in this fragment. Then we call this execution fragment a *non-blocking fragment* for  $R_i$  at  $s_j$  (Fig. 1a).
- When the context is clear, we omit the first subscript of  $F$  that denotes the reader. For instance, for a READ transaction  $R$ ,  $F_x(\alpha)^{(x_0)}$  denotes the non-blocking fragment of  $R$  on  $s_x$ .
2. Suppose READ transaction  $R_i$  completes in  $\alpha$ . Consider the execution fragment in  $\alpha$  between the event  $INV(R_i)$  and whichever of the events  $send(m_y^{r_i})_{r_i, s_y}$  and  $send(m_x^{r_i})_{r_i, s_x}$  that occurs later. If all the actions in this fragment occur at  $r_i$ , then we denote this fragment as  $I_i(\alpha)$  (Fig. 1a).
  3. Suppose READ transaction  $R_i$  completes in  $\alpha$ . Consider the execution fragments in  $\alpha$  that occurs between the later of the events  $recv(x)_{s_x, r_i}$  or  $recv(y)_{s_y, r_i}$ , i.e., at the point in  $\alpha$  when  $r_i$  receives responses from both servers, and the event  $RESP(R_i)$ . If all the actions in this fragment occur at  $r_i$ , then we denote this fragment by  $E_i(\alpha)^{(x, y)}$ , where  $R_i$  returns the values  $(x, y)$  (Fig. 1a) to the external client.
  4. We use  $R(\alpha)$  and  $W(\alpha)$  to denote the READ and WRITE transactions in the context of  $\alpha$ . When the the context is clear, we simply use  $R$  and  $W$ .
  5. We use the subscript of a returned value to denote the version identifier, which uniquely identifies a version from a totally ordered set. For instance,  $x_0$  is the  $0^{th}$  version of  $x$  (the initial value of object  $o_x$ ) on server  $s_x$ .

To play non-blocking fragments against the requirements of strict serializability and write isolation, we derive the following lemmas, which are frequently used by our proofs. In these lemmas, all READ transactions are assumed to have all SNOW properties, and we will derive contradictions in our main proofs. Due to space constraints, we explain these lemmas at a high level.

The following lemma proves that a READ transaction has to return the same version from both servers in order to satisfy strict serializability and write isolation.

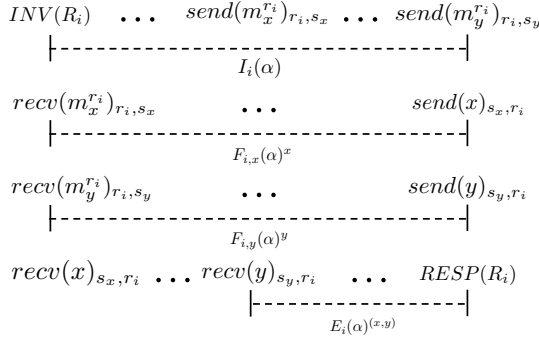
► **Lemma 3.** *Suppose  $\alpha$  is any execution of  $\mathcal{A}$  such that READ transaction  $R$  is in  $\alpha$ . Suppose the execution fragment  $I(\alpha) \circ F_x(\alpha)^{(x_t)} \circ F_y(\alpha)^{(y_s)} \circ E(\alpha)^{(x_{t'}, y_{s'})}$  in  $\alpha$ , corresponds to  $R$ , where  $x_t, x_{t'} \in V_1$  and  $y_s, y_{s'} \in V_2$ , then  $s = s' = t = t'$ .*

The following lemma states that we can create a new execution  $\alpha'$  that is indistinguishable to  $\alpha$  by swapping two fragments, which happen on two distinct automata in  $\alpha$  if either (a) both fragments have no input actions or (b) one of the fragments have no external (input or output) actions. Our proofs leverage this lemma to create new executions by swapping such fragments and finally derive an execution that violates strict serializability.

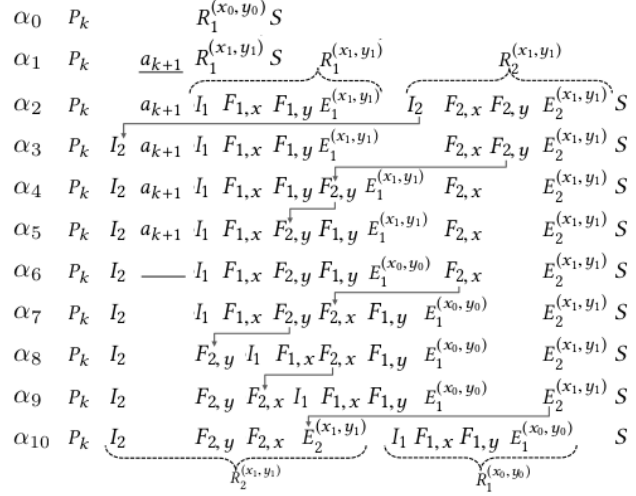
► **Lemma 4 (Commuting fragments).** *Let  $\alpha$  be an execution of  $\mathcal{A}$ . Suppose  $G_1(\alpha)$  and  $G_2(\alpha)$  are any execution fragments in  $\alpha$  such that all actions in each fragment occur only at one automaton and either (a) none of the fragments contain input actions, or (b) at least one of the fragments have no external actions. Suppose  $G_1(\alpha)$  and  $G_2(\alpha)$  occur at two distinct automata and the execution fragment  $G_1(\alpha) \circ G_2(\alpha)$  occurs in  $\alpha$ . Then there exists an execution  $\alpha'$  of  $\mathcal{A}$ , where the execution fragment  $G_2(\alpha) \circ G_1(\alpha)$  appears in  $\alpha'$ , such that (i)  $G_1(\alpha) \sim G_1(\alpha')$  and  $G_2(\alpha) \sim G_2(\alpha')$  (ii) the prefix in  $\alpha$  before  $G_1(\alpha) \circ G_2(\alpha)$  is identical to the prefix in  $\alpha'$  before  $G_1(\alpha') \circ G_2(\alpha')$ ; and (ii) the suffix in  $\alpha$  after  $G_1(\alpha) \circ G_2(\alpha)$  is identical to the suffix in  $\alpha'$  after the execution fragment  $G_2(\alpha') \circ G_1(\alpha')$ .*

The following lemma states that if there are two fair executions of  $\mathcal{A}$  with READ transaction  $R$  in each of them, and suppose at any server the non-blocking fragments of  $R$





(a) The relevant actions in the execution fragments  $I_i(\alpha)$ ,  $F_{i,x}(\alpha)^x$ ,  $F_{i,y}(\alpha)^y$  and  $E_i(\alpha)^{(x,y)}$  for any READ transaction  $R_i$ ,  $i \in \{1, 2\}$  of a fair execution  $\alpha$  of  $\mathcal{A}$ .



(b) Executions of  $\mathcal{A}$  with three clients and operations  $W$ ,  $R_1$  and  $R_2$  leading to the contradiction of  $S$  in  $\alpha_{10}$ . Arrows show the transposition of execution fragments from the previous execution.

are identical (in terms of the sequence of states and actions), then  $R$  returns the same value in both executions.

► **Lemma 5 (Indistinguishability).** *Let  $\alpha$  and  $\beta$  be executions of  $\mathcal{A}$  and let  $R$  be any READ transaction. Then (i) if  $F_x(\alpha) \stackrel{s_x}{\sim} F_x(\beta)$  then both  $R(\alpha)$  and  $R(\beta)$  respond with the same value  $x$  at  $s_x$ ; and (ii) if  $F_y(\alpha) \stackrel{s_y}{\sim} F_y(\beta)$  then both  $R(\alpha)$  and  $R(\beta)$  respond with the same value  $y$  at  $s_y$ .*

The following lemma shows that for any finite execution of  $\mathcal{A}$  that ends with the invocation of READ transaction  $R_1$ , it is always possible to have an execution of  $\mathcal{A}$  where the fragments  $I$ ,  $F_x$ ,  $F_y$  and  $E$  appear consecutively due to the asynchronous network.

► **Lemma 6.** *If any finite execution of  $\mathcal{A}$  ends with  $INV(R)$ , for a READ transaction  $R_1$  then there exists an extension  $\alpha$  which is a fair execution of  $\mathcal{A}$  and is of the form  $P(\alpha) \circ I(\alpha) \circ F_{1,x}(\alpha)^{(x)} \circ F_{1,y}(\alpha)^{(y)} \circ E(\alpha)^{(x,y)} \circ S(\alpha)$ , where  $P(\alpha)$  is the prefix and  $S(\alpha)$  denotes the rest of the execution.*

## 4 No SNOW with Three Clients and C2C

This section provides a formal proof of the SNOW Theorem with 3 clients, i.e., SNOW is impossible in a system with 3 or more clients even when client-to-client communication is allowed. The main result of this section is captured by the following theorem.

► **Theorem 7.** *The SNOW properties cannot be implemented in a system with two readers and one writer, for two servers even in the presence of client-to-client communication.*

Our proof strategy is to assume the existence of an algorithm  $\mathcal{A}$  that satisfies all SNOW properties and create an execution  $\alpha$  of  $\mathcal{A}$  that contradicts the S property. We begin with an execution of  $\mathcal{A}$  that contains READ transactions  $R_1$  and  $R_2$ , which both read  $s_x$  and  $s_y$ , and WRITE transaction  $W$  that writes  $(x_1, y_1)$  to  $s_x$  and  $s_y$  respectively (both servers

have initial values  $x_0, y_0$ ).  $R_1$  begins after  $W$  completes, and  $R_2$  begins after  $R_1$  completes. By the S property both  $R_1$  and  $R_2$  should return  $(x_1, y_1)$ . Then we create a sequence of executions of  $\mathcal{A}$  (Fig. 1b), where we interchange the fragments until we finally reach an execution in which  $R_2$  completes before  $R_1$  begins, but  $R_2$  returns  $(x_1, y_1)$  and  $R_1$  returns  $(x_0, y_0)$  which contradicts the S property.

The following lemma shows that in an execution of  $\mathcal{A}$  with a WRITE transaction  $W$  and a READ transaction  $R_1$ , there exists a point in the execution such that if  $R_1$  is invoked before that point then  $R_1$  returns  $(x_0, y_0)$  and if  $R_1$  is invoked after that point then  $R_1$  returns  $(x_1, y_1)$ .

► **Lemma 8** (Existence of  $\alpha_0$  and  $\alpha_1$ ). *There exist executions  $\alpha_0$  and  $\alpha_1$  of  $\mathcal{A}$  that contain transactions  $W$  and  $R_1$  that satisfy the following properties where  $k$  is some positive integer and  $a_1, \dots, a_k$  is a prefix of  $a_1, \dots, a_{k+1}$ :*

- (i)  $\alpha_0$  can be written as  $a_1, \dots, a_k \circ R_1(\alpha_0)^{(x_0, y_0)} \circ S(\alpha_0)$  ;
- (ii)  $\alpha_1$  can be written as  $a_1, \dots, a_{k+1} \circ R_1(\alpha_1)^{(x_1, y_1)} \circ S(\alpha_1)$ ; and
- (iii)  $a_{k+1}$  in  $\alpha_1$  occurs at  $r_1$ ,

**Proof.** Now we describe the construction of a sequence of finite executions of  $\mathcal{A}$ ,  $\{\gamma_k\}_{k=0}^\infty$  such that each  $\gamma_k$  contains  $W$  and  $R_1$ . Consider an execution  $\alpha$  of  $\mathcal{A}$  that contains  $W$ . Suppose  $R_1$  is invoked at  $r_1$  after the execution fragment  $a_1, \dots, a_{k+1}$ , a prefix of  $\alpha$ . Allowed by network asynchrony, let  $INV(R)$  be followed by only internal and external actions at  $r_1$  until both  $send(m_x^{r_1})_{r_1, s_x}$  and  $send(m_y^{r_1})_{r_1, s_y}$  occur, thereby creating an execution fragment of the form  $a_1, \dots, a_{k+1} \circ I_1(\alpha)$ . We denote  $a_1, \dots, a_{k+1}$  by  $P_{k+1}$ .

Next, consider the network delivers the message  $m_x^{r_1}$  at  $s_x$ , and delays all actions at other automata and also any input action at  $s_x$  until  $s_x$  sends  $x$  to  $r_1$ . Therefore, we achieve the execution fragment  $P_{k+1} \circ I_{1,x}(\alpha) \circ F_{1,x}(\alpha)$  of  $\mathcal{A}$ . Next, the network delivers  $m_y^{r_1}$  at  $s_y$  and delays all actions at other automata and input actions at  $s_y$  until  $s_y$  sends  $y$  to  $r_1$ . Then the network delivers  $x$  and  $y$  at  $r_1$  but it delays actions at other automata and any other input action at  $r_1$  until  $RESP(R_1)$  occurs. Now we have an execution fragment of  $\mathcal{A}$ , which can be written as  $P_{k+1} \circ I_1(\alpha) \circ F_{1,x}(\alpha)^{(x)} \circ F_{1,y}(\alpha)^{(y)} \circ E_1(\alpha)^{(x,y)}$ , where  $R_1$  responds with  $(x, y)$  such that  $(x, y) \in \{(x_0, y_0), (x_1, y_1)\}$ . We denote this finite execution prefix as  $\gamma_k$ . Therefore, there exists a sequence of such finite executions  $\{\gamma_k\}_{k=0}^\infty$ .

Because  $R_1$  precedes  $W$ , by the S property  $R_1$  must respond with  $(x_0, y_0)$  in  $\gamma_0$ . If  $k$  is large enough such that  $a_k$  occurs in  $\alpha$  after the completion of  $W$  then by the S property,  $R_1$  must return  $(x_1, y_1)$  in  $\gamma_{k+1}$  due to the S property. Therefore, there exists a minimum  $k$  where in  $\gamma_k$  READ transaction  $R_1$  returns  $(x_0, y_0)$  and in  $\gamma_{k+1}$ ,  $R_1$  returns  $(x_1, y_1)$ . Note that  $\gamma_k$  corresponds to  $\alpha_0$  and  $\gamma_{k+1}$  corresponds to  $\alpha_1$  in (i) and (ii) respectively.

Now, we prove case (iii) by eliminating the possibility of  $a_{k+1}$  occurring at  $s_x, s_y, w$  or  $r_2$ . The S property requires that  $R_1$  must return the same version from both  $s_x$  and  $s_y$ , which implies that  $s_x$  and  $s_y$  must send values of the same version. Observe that  $R_1$  returns the 0<sup>th</sup> version in  $\alpha_0$  and the 1<sup>st</sup> version in  $\alpha_1$ , while the prefixes  $P_k$  and  $P_{k+1}$  differ by a single action  $a_{k+1}$ . Importantly, just one action at any of  $s_x, s_y, r_2$  or  $w$  is not enough for  $s_x$  and  $s_y$  to coordinate the same version to send. Therefore,  $a_{k+1}$  must occur at  $r_1$ , which can possibly help coordinate by sending some information via  $m_x$  and  $m_y$  sent to  $s_x$  and  $s_y$  respectively.

*Case  $a_{k+1}$  occurs at  $s_x$ :* Consider the prefix of execution  $\alpha_0$  up to  $a_k$ . Suppose the network invokes  $R_1$  immediately after action  $a_k$  via  $INV(R_1)$ . By Lemma 6 there exists an execution  $\alpha'$  that contains an execution fragment of the form  $P_k \circ I_1(\alpha') \circ F_{1,x}(\alpha')^{(x)} \circ F_{1,y}(\alpha')^{(y)} \circ E(\alpha')^{(x,y)}$ . Then,  $I_1(\alpha_1) \stackrel{r_1}{\sim} I_1(\alpha')$  and  $F_{1,y}(\alpha_1) \stackrel{s_y}{\sim} F_{1,y}(\alpha')$  because in both



executions the actions of  $I_1$  occur entirely at  $r_1$  and those of  $F_{1,y}$  occur entirely at  $s_y$ , and thus they are unaffected by the addition of the single action  $a_{k+1}$  at  $s_x$ . As a result,  $F_{1,y}(\alpha')$  must send the same value  $y_1$  to  $r_1$  as in  $F_{1,y}(\alpha_1)$ . Then in  $\alpha'$ ,  $R_1(\alpha')$  returns  $y_1$  by Lemma 5, and thus  $R_1(\alpha')$  returns  $(x_1, y_1)$  by the S property. However, this contradicts the definition of  $k$ , i.e., the minimum value of  $k$  such that  $R_1$  responds with  $(x_0, y_0)$ .

Case  $a_{k+1}$  occurs at  $s_y$ : A contradiction can be shown by following a line of reasoning similar to the preceding case.

Case  $a_{k+1}$  occurs at  $w$ : This can be argued in a similar manner as the previous case with the trivial fact that  $F_{1,x}(\alpha_1) \stackrel{s_x}{\sim} F_{1,x}(\alpha')$  and  $F_{1,y}(\alpha_1) \stackrel{s_y}{\sim} F_{1,y}(\alpha')$ .

Case  $a_{k+1}$  occurs at  $r_2$ : A contradiction can be derived using a line of reasoning as in the previous case.

So we conclude that  $a_{k+1}$  must occur at  $r_1$  in  $\alpha_1$ . ◀

In the remainder of the section, we suppress the explicit reference to the execution. For instance, we use  $I_i$ ,  $F_{i,x}^{(x)}$ ,  $F_{i,y}^{(y)}$ ,  $E_i^{(x,y)}$  and  $S$  instead of  $I_i(\alpha)$ ,  $F_{i,x}(\alpha)^{(x)}$ ,  $F_{i,y}(\alpha)^{(y)}$ ,  $E_i(\alpha)^{(x,y)}$  and  $S(\alpha)$ . If a READ transaction  $R_i$  has an execution fragment of the form  $I_i \circ F_{i,x}^{(x)} \circ F_{i,y}^{(y)} \circ E_i^{(x,y)}$  we denote it as  $R_i^{(x,y)}$ . In the rest of the section,  $\alpha_0$ ,  $\alpha_1$ , and the value of  $k$  are the same as in the discussion above. We denote the execution fragments  $a_1, \dots, a_k$  and  $a_1, \dots, a_{k+1}$  as  $P_k$  and  $P_{k+1}$  respectively.

Our proof proceeds with a set of lemmas. Due to space limit, we omit the proof of some lemmas from the main paper, and present them in Appendix D as they are relatively straightforward. The first lemma proves there exists an execution in which two consecutive READ transactions follow a WRITE transaction, and both READ transactions return the new values by the WRITE transaction.

► **Lemma 9** (Existence of  $\alpha_2$ ). *There exists an execution  $\alpha_2$  of  $\mathcal{A}$  that contains  $W$ ,  $R_1$ , and  $R_2$ , and can be written in the form  $P_{k+1} \circ R_1^{(x_1, y_1)} \circ R_2^{(x_1, y_1)} \circ S$ , where both  $R_1$  and  $R_2$  return  $(x_1, y_1)$ .*

Based on the previous execution, the following lemma proves that there is an execution of  $\mathcal{A}$  where  $I_2$  occurs earlier than the action  $a_{k+1}$  and the invocation of  $R_1$ .

► **Lemma 10** (Existence of  $\alpha_3$ ). *There exists execution  $\alpha_3$  of  $\mathcal{A}$  that contains transactions  $W$ ,  $R_1$  and  $R_2$ , and can be written in the form  $P_k \circ I_2 \circ a_{k+1} \circ R_1^{(x_1, y_1)} \circ F_{2,x} \circ F_{2,y} \circ E_2 \circ S$ , where both  $R_1$  and  $R_2$  return  $(x_1, y_1)$ .*

**Proof.** Consider the execution  $\alpha_2$  as in Lemma 9. In the execution fragment  $I_1 \circ F_{1,x}^{(x_1)} \circ F_{1,y}^{(y_1)} \circ E_1^{(x_1, y_1)}$  in  $\alpha_2$ , none of the actions occur at  $r_2$  and by Lemma 8,  $a_{k+1}$  occurs at  $r_1$ , also the actions in  $I_2$  occur only at  $r_2$ . Starting with  $\alpha_2$ , and by repeatedly using Lemma 4, we create a sequence of four executions of  $\mathcal{A}$  by repeatedly swapping  $I_2$  with the execution fragments  $E_1^{(x_1, y_1)}$ ,  $F_{1,y}^{(y_1)}$ ,  $F_{1,x}^{(x_1)}$  and  $I_1$ , which appears in  $I_1 \circ F_{1,x}^{(x_1)} \circ F_{1,y}^{(y_1)} \circ E_1^{(x_1, y_1)} \circ I_2$ , where the following sequence of execution fragments  $I_1 \circ F_{1,x}^{(x_1)} \circ F_{1,y}^{(y_1)} \circ I_2 \circ E_1^{(x_1, y_1)}$  (by commuting  $I_2$  and  $E_1^{(x_1, y_1)}$ );  $I_1 \circ F_{1,x}^{(x_1)} \circ I_2 \circ F_{1,y}^{(y_1)} \circ E_1^{(x_1, y_1)}$  (by commuting  $I_2$  and  $F_{1,y}^{(y_1)}$ );  $I_1 \circ I_2 \circ F_{1,x}^{(x_1)} \circ F_{1,y}^{(y_1)} \circ E_1^{(x_1, y_1)}$  (by commuting  $I_2$  and  $F_{1,x}^{(x_1)}$ ) appear. Finally, we have an execution  $\alpha'$  of the form  $P_{k+1} \circ I_2 \circ R_1^{(x_1, y_1)} \circ F_{2,x}^{(x_1)} \circ F_{2,y}^{(y_1)} \circ E_2^{(x_1, y_1)} \circ S$  (by commuting  $I_2$  and  $I_1$ ) Next, from  $\alpha'$ , by using Lemma 4 and swapping  $a_{k+1}$  with  $I_2$  we have shown the existence of an execution  $\alpha_3$ . ◀

In the following lemma, we show that we can create an execution  $\alpha_4$  of  $\mathcal{A}$ , where  $F_{2,y}$  occurs immediately before  $E_1^{(x_1, y_1)}$ , while  $R_1$  and  $R_2$  both return  $(x_1, y_1)$ .

386 ► **Lemma 11** (Existence of  $\alpha_4$ ). *There exists execution  $\alpha_4$  of  $\mathcal{A}$  that contains transactions  $W$ ,  
 387  $R_1$  and  $R_2$  and can be written in the form  $P_k \circ I_2 \circ a_{k+1} \circ I_1 \circ F_{1,x} \circ F_{1,y} \circ F_{2,y} \circ E_1 \circ F_{2,x} \circ E_2 \circ S$ ,  
 388 where both  $R_1$  and  $R_2$  return  $(x_1, y_1)$ .*

389 Next, we create an execution  $\alpha_5$  where  $F_{2,y}$  occurs before  $F_{1,y}$ .

390 ► **Lemma 12** (Existence of  $\alpha_5$ ). *There exists execution  $\alpha_5$  of  $\mathcal{A}$  that contains transactions  $W$ ,  
 391  $R_1$  and  $R_2$  and can be written in the form  $P_k \circ I_2 \circ a_{k+1} \circ I_1 \circ F_{1,x} \circ F_{2,y} \circ F_{1,y} \circ E_1 \circ F_{2,x} \circ E_2 \circ S$ ,  
 392 where both  $R_1$  and  $R_2$  return  $(x_1, y_1)$ .*

393 **Proof.** Given all actions in  $F_{1,y}$  and  $F_{2,y}$  occur at  $s_y$  in  $\alpha_4$ , consider the prefix of  $\alpha_4$  that  
 394 ends with  $F_{1,x}$ . We extend this prefix as follows. In this prefix, the actions  $send(m_y^{r_2})_{r_2, s_y}$   
 395 and  $send(m_y^{r_1})_{r_1, s_y}$  do not have their corresponding *recv* actions. Suppose the network  
 396 delivers  $m_y^{r_2}$  at  $s_y$  (via the action  $recv(m_y^{r_2})_{r_2, s_y}$ ) and delays all actions, other than internal  
 397 and output actions at  $s_y$ , until  $s_y$  responds with  $y$ , via action  $send(y)_{s_y, r_2}$ . This extended  
 398 execution fragment is of the form  $F_{2,y}$ . Similarly, the network further extends the execution  
 399 by placing the action  $recv(m_y^{r_1})_{r_1, s_y}$  at  $s_y$  and creates the execution fragment of the form  
 400  $F_{1,y}$ . Note that, so far, the actions due to the above extensions are entirely at  $s_y$ . Suppose  
 401 the network makes the execution fragments  $E_1$  happen next by delivering values sent during  
 402  $F_{1,x}$  and  $F_{1,y}$  via the actions  $recv(x)_{s_x, r_1}$  and  $recv(y)_{s_y, r_1}$  respectively at  $r_1$ . Then  $F_{2,x}$   
 403 occurs next, such that this fragment contains exactly the same sequence of actions as in the  
 404 corresponding execution fragment in  $\alpha_4$ . This is possible because they are not influenced by  
 405 any output action in  $F_{2,y}$  or  $F_{1,y}$ . Suppose the network places the execution fragment  $E_2$   
 406 next. Let us denote the execution that is an extension of this finite execution so far as  $\alpha_5$ ,  
 407 which is of the form  $P_k \circ I_2 \circ a_{k+1} \circ I_1 \circ F_{1,x} \circ F_{2,y} \circ F_{1,y} \circ E_1 \circ F_{2,x} \circ E_2 \circ S$ . Now we need to  
 408 argue about the values returned by the reads.

409 Note that both  $\alpha_4$  and  $\alpha_5$  have the same execution fragment  $F_{1,x}(\alpha_4)$ . Therefore,  
 410  $F_{1,x}(\alpha_4) \stackrel{s_x}{\sim} F_{1,x}(\alpha_5)$ , and thus  $s_x$  also returns  $x_1$  in  $F_{1,x}$  in  $\alpha_5$ . Next by Lemma 5 for  $R_1$ ,  $s_y$   
 411 returns  $y_1$  in  $F_{1,y}$  and hence by the S property,  $R_1(\alpha_5)$  returns  $(x_1, y_1)$ , i.e., that  $r_1$  returns  
 412 the new version of object values. Therefore,  $F_{1,x}(\alpha_4)$ ,  $F_{1,y}$  and  $E_1$  are of the form  $F_{1,x}^{(x_1)}$ ,  
 413  $F_{1,y}^{(y_1)}$  and  $E_1^{(x_1, y_1)}$ , respectively.

414 Note that by construction of  $\alpha_5$  above, the execution fragment  $F_{2,x}$  in both  $\alpha_4$  and  $\alpha_5$   
 415 is the same, therefore,  $F_{2,x}(\alpha_4) \stackrel{s_x}{\sim} F_{2,x}(\alpha_5)$ . Hence as in  $\alpha_4$ ,  $s_x$  returns  $x_1$  in the execution  
 416 fragment  $F_{2,x}(\alpha_5)$  in  $\alpha_5$ , i.e., of the form  $F_{2,x}(\alpha_5)^{(x_1)}$ . Since  $s_x$  returns  $x_1$  in  $F_{2,x}$  in  $\alpha_5$ , by  
 417 Lemma 5 and the S property,  $R_2$  returns  $(x_1, y_1)$  and hence  $E_2$  is of the form  $E_2^{(x_1, y_1)}$ .

418 From the above argument we know that  $\alpha_5$  is of the form  $P_k \circ I_2 \circ a_{k+1} \circ I_1 \circ F_{1,x}^{(x_1)} \circ F_{2,y}^{(y_1)} \circ$   
 419  $F_{1,y}^{(y_1)} \circ E_1^{(x_1, y_1)} \circ F_{2,x}^{(x_1)} \circ E_2^{(x_1, y_1)} \circ S$ . ◀

420 In the next lemma, we show the existence of an execution of  $\mathcal{A}$  where  $R_1$  returns  $(x_0, y_0)$   
 421 and  $I_2$  occurs immediately after  $a_k$  and  $R_2$  responds with  $(x_1, y_1)$ .

422 ► **Lemma 13** (Existence of  $\alpha_6$ ). *There exists execution  $\alpha_6$  of  $\mathcal{A}$  that contains transactions  $W$ ,  
 423  $R_1$  and  $R_2$  and can be written in the form  $P_k \circ I_2 \circ I_1 \circ F_{1,x} \circ F_{2,y} \circ F_{1,y} \circ E_1 \circ F_{2,x} \circ E_2 \circ S$ ,  
 424 where  $R_1$  returns  $(x_0, y_0)$  and  $R_2$  returns  $(x_1, y_1)$ .*

425 **Proof.** The crucial part of this proof is to carefully use the result of Lemma 8 so that  $R_1$   
 426 returns  $(x_0, y_0)$ , instead of  $(x_1, y_1)$ . Note that the same prefix  $P_k$  appears in  $\alpha_5$  of Lemma 12  
 427 as well as in  $\alpha_0$  and  $\alpha_1$  of Lemma 8, where  $k$  is defined as in Lemma 8.

428 By Lemma 8, action  $a_{k+1}$  occurs at  $r_1$ . In the execution fragment  $a_{k+1} \circ I_1 \circ F_{1,x}^{(x_1)} \circ F_{2,y}^{(y_1)}$   
 429 of  $\alpha_5$ , the actions in  $a_{k+1} \circ I_1$  occur at  $r_1$ ; actions in  $F_{1,x}^{(x_1)}$  occur at  $s_x$ ; and actions in  $F_{2,y}^{(y_1)}$

occur at  $s_y$ . Now consider the prefix of execution  $\alpha_4$  ending with  $I_2$  and the network invokes  $R_1$  immediately after  $I_2$  (instead of after  $a_{k+1}$ ) and extends it by the execution fragment  $I_1 \circ F_{1,x} \circ F_{2,y}$  to create a new finite execution  $\epsilon$ , which is of the form  $P_k \circ I_2 \circ I_1 \circ F_{1,x} \circ F_{2,y}$ . As a result,  $a_{k+1}$  may not be in  $\epsilon$  because we introduce changes before  $a_{k+1}$  occurs.

Note that if in the prefix  $P_k \circ I_2(\epsilon) \circ I_1(\epsilon) \circ F_{1,x}(\epsilon) \circ F_{2,y}(\epsilon)$  of  $\epsilon$  we ignore the actions in  $I_2(\epsilon)$  then the remaining execution is the same as the prefix  $P_k \circ I_1(\alpha_0) \circ F_{1,x}(\alpha_0) \circ F_{2,y}(\alpha_0)$  of  $\alpha_0$  in Lemma 8. Here we explicitly use the notations  $\epsilon$  and  $\alpha_0$  to avoid confusion. Since the actions in  $I_2(\epsilon)$  have no impact on the actions in  $I_1(\epsilon) \circ F_{1,x}(\epsilon) \circ F_{2,y}(\epsilon)$ , we have  $F_{1,x}(\epsilon) \stackrel{s_x}{\sim} F_{1,x}(\alpha_0)$ . Therefore, by Lemma 3  $F_{1,x}(\epsilon)$  returns  $x_0$  as in  $F_{1,x}(\alpha_0)$ , i.e.,  $s_x$  returns  $x_0$  in  $F_{1,x}$ . Now by Lemma 5 we conclude that for any extension of  $\epsilon$ , say  $\gamma$ , READ transaction  $R_1(\gamma)$  returns  $x_0$  at  $s_x$  and by the S property  $R_1(\gamma)$  returns  $(x_0, y_0)$ . Also, since  $F_{2,y}(\alpha_5) \stackrel{s_y}{\sim} F_{2,y}(\epsilon) \stackrel{s_y}{\sim} F_{2,y}(\gamma)$  by Lemma 5 and the S property,  $R_2(\gamma)$  must return  $(x_1, y_1)$ . Therefore,  $\gamma$  has an extension of  $\alpha_6$  (Fig. 1b) which is of the form  $P_k \circ I_2 \circ I_1 \circ F_{1,x}^{(x_0)} \circ F_{2,y}^{(y_1)} \circ F_{1,y}^{(y_0)} \circ E_1^{(x_0, y_0)} \circ F_{2,x}^{(x_1)} \circ E_2^{(x_1, y_1)} \circ S$  as in the statement of the lemma.  $\blacktriangleleft$

The following lemma shows that there exists an execution  $\alpha_7$  for  $\mathcal{A}$  where  $F_{2,x}$  appears before  $F_{1,y} \circ E_1$ , where  $R_1$  returns  $(x_0, y_0)$  and  $R_2$  returns  $(x_1, y_1)$ . The lemma can be proven by starting from  $\alpha_6$  in Lemma 13 and moving the execution fragments of  $R_2$  earlier, a little at a time, until finally we have  $R_2$  finishing before  $R_1$  starts. This simply uses commutativity since the actions in the swapped execution fragments occur at different automata.

► **Lemma 14** (Existence of  $\alpha_7$ ). *There exists execution  $\alpha_7$  of  $\mathcal{A}$  that contains transactions  $W$ ,  $R_1$  and  $R_2$ , and can be written in the form  $P_k \circ I_2 \circ I_1 \circ F_{1,x} \circ F_{2,y} \circ F_{2,x} \circ F_{1,y} \circ E_1 \circ E_2 \circ S$  where  $R_1$  returns  $(x_0, y_0)$  and  $R_2$  returns  $(x_1, y_1)$ .*

The following lemma leverages Lemma 4 to show the existence of an execution  $\alpha_8$  of  $\mathcal{A}$  where  $F_{2,y}$  appears before  $I_1 \circ F_{1,x}$ , and  $R_1$  returns  $(x_0, y_0)$  while  $R_2$  returns  $(x_1, y_1)$ .

► **Lemma 15** (Existence of  $\alpha_8$ ). *There exists execution  $\alpha_8$  of  $\mathcal{A}$  that contains transactions  $W$ ,  $R_1$  and  $R_2$  and can be written in the form  $P_k \circ I_2 \circ F_{2,y} \circ I_1 \circ F_{1,x} \circ F_{2,x} \circ F_{1,y} \circ E_1 \circ E_2 \circ S$ , where  $R_1$  returns  $(x_0, y_0)$  and  $R_2$  returns  $(x_1, y_1)$ .*

The following lemma shows the existence of an execution  $\alpha_9$ , of  $\mathcal{A}$ , where  $F_{2,x}$  appears before  $F_{1,x}$ .

► **Lemma 16** (Existence of  $\alpha_9$ ). *There exists execution  $\alpha_9$  of  $\mathcal{A}$  that contains transactions  $W$ ,  $R_1$  and  $R_2$  and can be written in the form  $P_k \circ I_2 \circ F_{2,y}^{(y_1)} \circ I_1 \circ F_{2,x}^{(x_1)} \circ F_{1,x}^{(x_0)} \circ F_{1,y}^{(y_0)} \circ E_1^{(x_0, y_0)} \circ E_2^{(x_1, y_1)} \circ S$  where  $R_1$  returns  $(x_0, y_0)$  and  $R_2$  returns  $(x_1, y_1)$ .*

Now we show the existence of an execution of  $\mathcal{A}$  where the execution fragments corresponding to  $R_2$  appears before  $R_1$ , where  $R_1$  returns  $(x_0, y_0)$  and  $R_2$  completes by returning  $(x_1, y_1)$ .

► **Lemma 17** (Existence of  $\alpha_{10}$ ). *There exists an execution  $\alpha_{10}$  of  $\mathcal{A}$  that contains transactions  $W$ ,  $R_1$  and  $R_2$  and can be written in the form  $P_k \circ R_2^{(x_1, y_1)} \circ R_1^{(x_0, y_0)} \circ S$ . where  $R_1$  returns  $(x_0, y_0)$  and  $R_2$  returns  $(x_1, y_1)$ .*

**Proof.** Now, by applying Lemma 4 to  $\alpha_9$ , we can swap  $F_{2,x}$  and  $I_1$  to create an execution  $\alpha_{10}$  (Fig. 1b) of  $\mathcal{A}$ , which is of the form  $P_k \circ I_2 \circ F_{2,y}^{(y_1)} \circ F_{2,x}^{(x_1)} \circ R_1^{(x_0, y_0)} \circ E_2^{(x_1, y_1)} \circ S$ , where the returned values are determined by Lemma 3.

Note that none of the actions in  $I_1 \circ F_{1,x}^{(x_0)} \circ F_{1,y}^{(y_0)} \circ E_1^{(x_0, y_0)}$  occur at  $r_2$  and all actions in  $E_2^{(x_1, y_1)}$  occur at  $r_2$ . Therefore, by applying Lemma 4, we can consecutively swap  $E_2$  with

473  $E_1, F_{1,y}, I_1$ , and  $F_{1,x}$ . Therefore, we create a sequence of four executions of  $\mathcal{A}$  to arrive at  
 474 execution  $\alpha_{10}$  (Fig. 1b) of the form  $P_k \circ R_2^{(x_1, y_1)} \circ R_1^{(x_0, y_0)} \circ S$ . ◀

475 Using the above results, we prove Theorem 7 for 3-clients by showing the existence of  
 476  $\alpha_{10}$ , where  $R_2$  completes before  $R_1$  is invoked and  $R_2$  returns  $(x_1, y_1)$  whereas  $R_1$  returns  
 477  $(x_0, y_0)$ , which violates the S property.

478 **Proof.**  $\alpha_{10}$  as shown in Fig. 1b provides a contradicting execution, where  $R_1$  returns  $(x_0, y_0)$   
 479 and  $R_2$  returns  $(x_1, y_1)$ , but  $R_1$  is in real time after  $R_2$ , which violates the S property. ◀

## 480 5 Two Client Open Question

481 This section closes the open question of whether SNOW properties can be implemented with  
 482 two clients. We first prove that SNOW remains impossible in a 2-client 2-server system if  
 483 the clients cannot directly send messages to each other. However, in the presence of client to  
 484 client communication, it is possible to have all SNOW properties with two clients and at  
 485 least two servers.

### 486 5.1 No SNOW Without C2C Messages

487 In this section, we prove the following results that states it is impossible guarantee the  
 488 SNOW properties in a transaction processing system with two clients, without client-to-client  
 489 communication.

490 ► **Theorem 18.** *The SNOW properties cannot be implemented in a system with two clients  
 491 and two servers, where the clients do not communicate with each other.*

492 We use the same system model as in Section 4: two servers  $s_x$  and  $s_y$  with two clients,  
 493 a reader  $r_1$  that issues only READ transactions and a writer  $w$  that issues only WRITE  
 494 transactions. A WRITE transaction  $W$  writes  $(x_1, y_1)$  to  $s_x$  and  $s_y$ , and a READ transaction  
 495  $R$  reads both servers. We assume that there is a bi-directional communication channel  
 496 between any pair of client and server and any pair of servers. There is no communication  
 497 channel between clients. We assume that each transaction can be identified by a unique  
 498 number, e.g., transaction identifier.

499 Our strategy is still proof by contradiction: We assume there exists some algorithm  $\mathcal{A}$  that  
 500 satisfies all SNOW properties, and then we show the existence of a sequence of executions of  
 501  $\mathcal{A}$ , eventually leading to an execution that contradicts the S property. First, we show the  
 502 existence of an execution  $\alpha$  of  $\mathcal{A}$  where  $R_1$  is invoked after  $W$  completes, where the send  
 503 actions  $send(m_x^{r_1})_{r_1, s_x}$  and  $send(m_y^{r_1})_{r_1, s_y}$  at the  $r_1$  occur consecutively in  $P(\alpha)$ , which is  
 504 a prefix of  $\alpha$ . Then we show that  $\alpha$  can be written in the form  $P(\alpha) \circ F_{1,x}(\alpha)$  (Fig. 2 (a),  
 505 Lemma 19). We then prove the existence of another execution  $\beta$ , which can be written in  
 506 the form  $P(\beta) \circ F_{1,x}(\beta) \circ F_{1,y}(\beta)$  by extending  $\alpha$  with an execution fragment  $F_{1,y}(\beta)$ , such  
 507 that  $F_{1,x}(\beta) \stackrel{s_x}{\sim} F_{1,x}(\alpha)$  (Fig. 2 (b); Lemma 20). Note that in any arbitrary extension of  $\beta$ ,  
 508  $R_1$  eventually returns  $(x_1, y_1)$ . Next, we show the existence of an execution  $\gamma$  of the form  
 509  $P(\gamma) \circ F_{1,x}(\gamma) \circ F_{1,y}(\gamma)$ , where the send actions  $send(m_x^{r_1})_{r_1, s_x}$  and  $send(m_y^{r_1})_{r_1, s_y}$  at  $r_1$  occur  
 510 before  $W$  is invoked (Fig. 2 (c);  $\gamma$ ), but  $F_{1,x}(\gamma)$  and  $F_{1,y}(\gamma)$  occur after  $RESP(W)$  as in  $\beta$ .  
 511 Based on  $\gamma$ , we show the existence of an execution  $\delta$  of the form  $P(\eta) \circ F_{1,x}(\eta) \circ F_{1,y}(\eta) \circ S(\eta)$ ,  
 512 where  $R_1$  responds with  $(x_1, y_1)$ . Finally, starting with  $\eta$ , we create a sequence of executions  
 513  $\delta(\equiv \eta), \delta^{(1)}, \dots, \delta^{(f)}$ , of  $\mathcal{A}$ , where in each of them  $R_1$  responds with  $(x_1, y_1)$  (Fig. ?? (e) and  
 514 (g); Lemma 18). Additionally, for any  $\delta^{(i)}$ , the fragments  $F_{1,x}(\delta^{(i)})$  and  $F_{1,y}(\delta^{(i)})$  appear

515 before  $\delta^{(i-1)}$ . Based on  $\delta^{(f)}$ , we prove the existence of an execution  $\phi$  (Fig. ?? (h)), where  
 516  $R_1$  returns  $(x_1, y_1)$  even before  $W$  begins, which violates the  $S$  property.

517 Now, we explain the relevant lemmas used by the main proof. Due to space constraints,  
 518 we present their proof in Appendix E.1. The first lemma states that there exists a finite  
 519 execution of  $\mathcal{A}$ , where  $R_1$  begins after  $W$  completes, and the actions  $send(m_x^{r_1})_{r_1, s_x}$  and  
 520  $send(m_y^{r_1})_{r_1, s_y}$  occur before either of the servers  $s_x$  and  $s_y$  receives  $m_x^{r_1}$  or  $m_y^{r_1}$  from  $r_1$ ; also,  
 521 server  $s_x$  responds to  $r_1$  in a non-blocking manner (execution  $\alpha$  in Fig. 2).

522 ► **Lemma 19.** *There exists a finite execution  $\alpha$  of  $\mathcal{A}$  that contains transactions  $R_1(\alpha)$  and  
 523  $W(\alpha)$  where  $INV(R)$  appears after  $RESP(W)$ , and the following conditions hold:*

- 524 (i) *The actions  $send(m_x^{r_1})_{r_1, s_x}$  and  $send(m_y^{r_1})_{r_1, s_y}$  appear consecutively in  $trace(\alpha)|_{r_1}$ ; and*  
 525 (ii)  *$\alpha$  contains the execution fragment  $F_{1,x}(\alpha)$ .*

526 The following lemma states that there is an execution  $\beta$  where  $R$  begins after  $W$  completes,  
 527 and the two send events at  $r_1$  occurs before  $F_{1,x}(\beta)$ , which thus occurs before  $F_{1,y}(\beta)$  ( $\beta$  in  
 528 Fig. 2).

529 ► **Lemma 20.** *There exists an execution  $\beta$  of  $\mathcal{A}$  that contains transactions  $R_1$  and  $W$  where  
 530  $INV(R_1)$  appears after  $RESP(W)$  and the following conditions hold:*

- 531 (i) *The actions  $send(m_x^{r_1})_{r_1, s_x}$  and  $send(m_y^{r_1})_{r_1, s_y}$  appear consecutively in  $trace(\beta)|_{r_1}$ ; and*  
 532 (ii)  *$\beta$  contains the execution fragment  $F_{1,x}(\beta) \circ F_{1,y}(\beta)$ .*

533 The following result shows that starting with  $\beta$  there is an execution  $\gamma$ , where  $R_1$  is  
 534 invoked before  $W$  is invoked while  $send(m_x^{r_1})_{r_1, s_x}$  and  $send(m_y^{r_1})_{r_1, s_y}$  occur before  $INV(W)$   
 535 (Fig. 2 (c)) and  $m_x^{r_1}$  and  $m_y^{r_1}$  from  $r_1$  reach  $s_x$  and  $s_y$  after the action  $RESP(W)$ .

536 ► **Lemma 21.** *There exists an execution  $\gamma$  that contains  $R_1$  and  $W$ , where the action  
 537  $INV(R_1)$  appears before  $INV(W)$  and  $RESP(R_1)$  appears after  $RESP(W)$ , and the following  
 538 conditions hold for  $\gamma$ :*

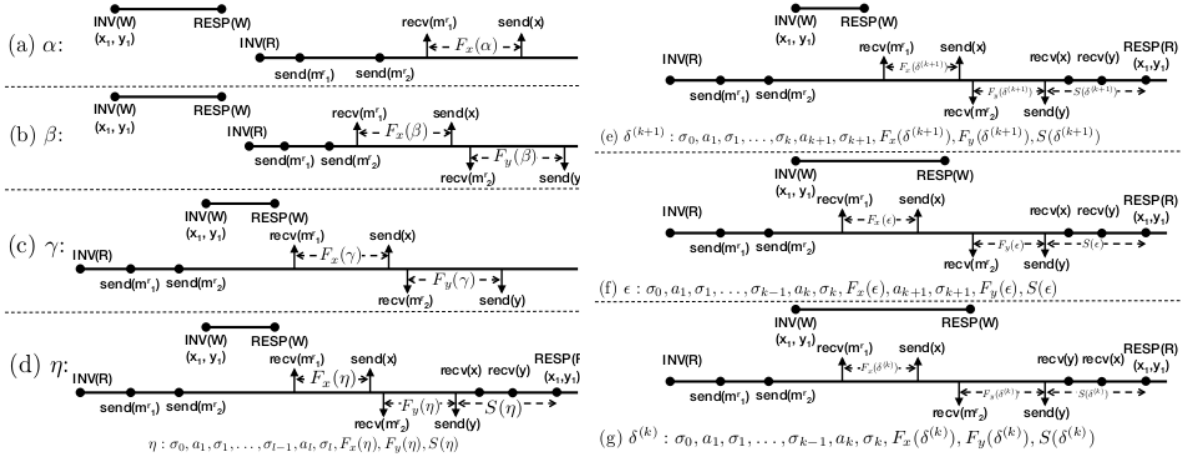
- 539 (i) *The actions  $send(m_x^{r_1})_{r_1, s_x}$  and  $send(m_y^{r_1})_{r_1, s_y}$  appear before  $INV(W)$  and they appear  
 540 consecutively in  $trace(\gamma)|_{r_1}$ ;*  
 541 (ii)  *$\gamma$  contains the execution fragment  $F_{1,x}(\gamma) \circ F_{1,y}(\gamma)$ ; and*  
 542 (iii) *action  $RESP(W)$  occurs before  $F_{1,x}(\gamma)$ .*

543 In the following lemma, we show there exists an execution  $\eta$  of  $\mathcal{A}$  of the form  $P(\eta) \circ$   
 544  $F_{1,x}(\eta) \circ F_{1,y}(\eta) \circ S(\eta)$  where  $RESP(R)$  appears in  $S(\eta)$  (Fig. 2 (d)) and  $R(\eta)$  returns  $(x_1, y_1)$ .

545 ► **Lemma 22.** *There exists an execution  $\eta$  of  $\mathcal{A}$  that contains transactions  $R$  and  $W$  where  
 546  $INV(R)$  appears before  $INV(W)$ ;  $RESP(R_1)$  appears after  $RESP(W)$  and the following  
 547 conditions hold for  $\eta$ :*

- 548 (i)  *$\eta$  can be written in the form  $P(\eta) \circ F_{1,x}(\eta) \circ F_{1,y}(\eta) \circ S(\eta)$ , for some  $P(\eta)$  and  $S(\eta)$ ;*  
 549 (ii) *The actions  $send(m_x^{r_1})_{r_1, s_x}$  and  $send(m_y^{r_1})_{r_1, s_y}$  appear before  $INV(W)$  and they appear  
 550 consecutively in  $trace(\eta)|_{r_1}$ ;*  
 551 (iii) *action  $RESP(W)$  occurs before  $F_{1,x}(\eta)$ ; and*  
 552 (iv)  *$R_1(\eta)$  returns  $(x_1, y_1)$ .*

553 The following proof proves Theorem 18. In the proof we start with an execution  $\eta$  and  
 554 create a sequence of executions of  $\mathcal{A}$ , where each one is of the form  $P(\cdot) \circ F_{1,x}(\cdot) \circ F_{1,y}(\cdot) \circ S(\cdot)$ ,  
 555 with progressively shorter  $P(\cdot)$  until we have a final execution that contradicts the  $S$  property.



■ **Figure 2** Schematic representation of executions  $\alpha$ ,  $\beta$ ,  $\gamma$  and  $\eta$ , of  $\mathcal{A}$  with transactions  $R$  and  $W$ . The executions evolve from left to right. The dots denote external events at clients. The up-arrow marks denote external actions at  $s_x$ . The down-arrow marks denote external actions at  $s_y$ .

**Proof.** Consider an execution  $\delta^{(\ell)}$  of  $\mathcal{A}$  as in Lemma 22, and let  $P(\delta^{(\ell)})$  be the execution fragment  $\sigma_0, a_1, \dots, a_\ell, \sigma_\ell$ . By Lemma 22,  $RESP(R_1(\delta^{(\ell)}))$  returns  $(x_1, y_1)$  and  $\delta^{(\ell)}$  is also of the form  $\sigma_0, a_1, \dots, a_\ell, \sigma_\ell \circ F_{1,x}(\delta^{(\ell)}) \circ F_{1,y}(\delta^{(\ell)}) \circ S(\delta^{(\ell)})$ .

Next, we inductively prove the existence of a finite sequence of executions of  $\mathcal{A}$ —i.e., by proving the existence of a new execution based on the existence of a previous one—as  $\delta^{(\ell)}, \delta^{(\ell-1)}, \dots, \delta^{(i)}, \delta^{(i-1)}, \dots, \delta^{(f)}$ , for some positive integer  $f$ , with the following properties: (a) Each of the execution in the sequence can be written in the form  $\sigma_0, a_1, \dots, a_i, \sigma_i \circ F_{1,x}(\delta^{(i)}) \circ F_{1,y}(\delta^{(i)}) \circ S(\delta^{(i)})$  (or  $P(\cdot) \circ F_{1,x}(\cdot) \circ F_{1,y}(\cdot) \circ S(\cdot)$ ); (b) for each  $i, f \leq i < \ell$ , we have  $P(\delta^{(i)})$  to be a prefix of  $P(\delta^{(i+1)})$ ; and (c)  $R_1(\delta^{(f)})$  returns  $(x_0, y_0)$ , and for any  $i, f < i \leq \ell$ , we have  $R_1(\delta^{(i)})$  that returns  $(x_1, y_1)$ . Note that there is a final execution of the form  $\delta^{(f)}$  because of the initial values of  $x_0$  and  $y_0$  and the WRITE transaction  $W$ .

Clearly, there exists an integer  $k, f \leq k < \ell$ , such that  $R(\delta^{(k)})$  returns  $(x_0, y_0)$  and  $R_1(\delta^{(k+1)})$  returns  $(x_1, y_1)$ . Now we start with execution  $\delta^{(k+1)}$  and construct an execution  $\delta^{(k)}$  as described in the rest of the proof. The following argument will show that  $R(\delta^{(k)})$  must also return  $(x_1, y_1)$ , which contradicts the assumption of having the  $S$  property.

Consider the execution  $\delta^{(k+1)}$  of the form  $\sigma_0, a_1, \dots, a_{k+1}, \sigma_{k+1} \circ F_{1,x}(\delta^{(k+1)}) \circ F_{1,y}(\delta^{(k+1)}) \circ S(\delta^{(k+1)})$ . The action  $a_{k+1}$  can occur at any of the automata  $r, w, s_x$ , and  $s_y$ . Therefore, we consider the following four possible cases.

Case (i)  $a_{k+1}$  occurs at  $w$ : The execution fragments  $F_{1,x}(\delta^{(k+1)})$  and  $F_{1,y}(\delta^{(k+1)})$  do not contain any input action at  $s_x$  or  $s_y$ .  $a_{k+1}$  does not occur at  $s_x$  or  $s_y$ . Therefore, the asynchronous network can delay the occurrence of  $a_{k+1}$  at  $w$  to create the finite execution  $\sigma_0, a_1, \dots, a_k, \sigma_k \circ F_{1,x}(\delta^{(k+1)}) \circ F_{1,y}(\delta^{(k+1)})$ , of. Also, there exists an execution  $\delta^{(k)}$  that is an extension of the above finite execution, where  $R_1$  completes in  $\delta^{(k)}$ . Clearly,  $\delta^{(k)}$  can be written as  $\sigma_0, a_1, \dots, a_k, \sigma_k \circ F_{1,x}(\delta^{(k)}) \circ F_{1,y}(\delta^{(k)}) \circ S(\delta^{(k)})$ , where  $S(\delta^{(k)})$  is the tail of the execution resulting from the extension. Moreover,  $F_{1,x}(\delta^{(k+1)})$  is indistinguishable from  $F_{1,x}(\delta^{(k)})$  at  $s_x$ , i.e.,  $F_{1,x}(\delta^{(k+1)}) \stackrel{s_x}{\sim} F_{1,x}(\delta^{(k)})$ . Therefore,  $send(x)_{s_x, r_1}$  has the same object value  $x$  in both fragments, which means  $R_1$  returns  $x_1$ , and thus  $R_1(\delta^{(k)})$  must return  $(x_1, y_1)$  by the property  $S$ .

Case (ii)  $a_{k+1}$  occurs at  $r$ : Similar to Case (i).

Case (iii)  $a_{k+1}$  occurs at  $s_x$ : Observe that the two execution fragments  $a_{k+1}\sigma_{k+1} \circ$



586  $F_{1,x}(\delta^{(k+1)})$  and  $F_{1,y}(\delta^{(k+1)})$  occur at separate automata, i.e., at  $s_x$  and  $s_y$  respectively.  
 587 Also, the execution fragments  $a_{k+1}\sigma_{k+1} \circ F_{1,x}(\delta^{(k+1)})$  and  $F_{1,y}(\delta^{(k+1)})$  do not contain any  
 588 input actions at  $s_x$  or  $s_y$ . Therefore, we can create an execution  $\epsilon$ , which can be expressed  
 589 as  $\sigma_0, a_1, \dots, a_k, \sigma_k \circ F_{1,y}(\epsilon) \circ a_{k+1}, \sigma_{k+1} \circ F_{1,x}(\epsilon) \circ S(\epsilon)$ . Clearly,  $F_{1,x}(\epsilon) \stackrel{s_x}{\sim} F_{1,x}(\delta^{(k+1)})$  and  
 590  $F_{1,y}(\epsilon) \stackrel{s_y}{\sim} F_{1,y}(\delta^{(k+1)})$ . Because  $send(x)_{s_x, r_1}$  occurs in both  $F_{1,x}(\epsilon)$  and  $F_{1,x}(\delta^{(k+1)})$ , it sends  
 591 the same value  $x_1$  to  $r_1$ . Therefore,  $R$  returns  $(x_1, y_1)$ .

592 Now let us denote the execution fragment  $\sigma_0, a_1, \dots, a_k, \sigma_k \circ F_{1,y}(\epsilon)$  by  $\epsilon'$ , which is simply  
 593 a finite prefix of  $\epsilon$ . Allowed by the asynchronous network, we append  $recv(m_x^{r_1})_{r_1, s_x}$  to  $\epsilon'$ ,  
 594 and create a finite execution  $\epsilon''$  as  $\sigma_0, a_1, \dots, a_k, \sigma_k \circ F_{1,y}(\epsilon''), recv(m_x^{r_1})_{r_1, s_x}$ , and delay any  
 595 input action at  $s_y$ .

596 Let  $\epsilon'''$  be an extension of  $\epsilon''$ . Clearly,  $F_{1,y}(\epsilon''') \stackrel{s_y}{\sim} F_{1,y}(\epsilon'')$ , and thus  $send(y)_{s_y, r_1}$  sends  
 597 the same value  $y_1$  in  $\epsilon''$  and  $\epsilon'''$ . By the  $N$  property,  $send(x)_{s_x, r_1}$  eventually occurs, and by  
 598 the  $O$  property  $x$  is sent to  $r_1$ . Therefore,  $R_1$  completes in  $\epsilon'''$ , which implies that  $R(\epsilon''')$   
 599 must return  $(x_1, y_1)$ .

600 Note that the execution fragment of  $\epsilon'''$  has no input actions of  $s_x$  between  $recv(m_x^{r_1})_{r_1, s_x}$   
 601 and  $send(x)_{s_x, r_1}$ , which can be identified as  $F_{1,x}(\epsilon''')$ . Therefore,  $\epsilon'''$  can be written as  
 602  $\sigma_0, a_1, \dots, a_k, \sigma_k \circ F_{1,y}(\epsilon''') \circ F_{1,x}(\epsilon''') \circ S(\epsilon''')$ .

603 Next, since  $F_{1,x}(\epsilon''')$  and  $F_{1,y}(\epsilon''')$  contain actions of different automata, we can create an  
 604 execution prefix  $\epsilon^{(iv)}$  as  $\sigma_0, a_1, \dots, a_k, \sigma_k \circ F_{1,x}(\epsilon^{(iv)}) \circ F_{1,y}(\epsilon^{(iv)})$ , where  $F_{1,x}(\epsilon^{(iv)})$  appears  
 605 before  $F_{1,y}(\epsilon^{(iv)})$ . Next, we create an execution  $\delta^{(k)}$  as an extension of  $\epsilon^{(iv)}$ , which can be  
 606 written as  $\sigma_0, a_1, \dots, a_k, \sigma_k \circ F_{1,x}(\delta^{(k)}) \circ F_{1,y}(\delta^{(k)}) \circ S(\delta^{(k)})$ . Then by the  $O$  and  $N$  properties,  
 607  $R_1$  completes in  $\delta^{(k)}$ . Since  $F_{1,x}(\epsilon^{(iv)}) \stackrel{s_x}{\sim} F_{1,x}(\epsilon''')$ ,  $send(x)_{s_x, r_1}$  returns  $x_1$  in  $\epsilon^{(iv)}$ . Similarly,  
 608 because  $F_{1,x}(\delta^{(k)}) \stackrel{s_x}{\sim} F_{1,x}(\epsilon^{(iv)})$ ,  $send(x)_{s_x, r_1}$  returns  $x_1$  in  $\delta^{(k)}$ . So,  $R_1(\delta^{(k)})$  returns  $(x_1, y_1)$ .

609 Case (iv)  $a_{k+1}$  occurs at  $s_y$ : Because  $a_{k+1}$  occurs at server  $s_y$  and  $F_{1,x}(\delta^{(k+1)})$  occurs  
 610 at server  $s_x$  (different automata), we can create a new execution  $\epsilon$  of  $\mathcal{A}$  (Fig. ?? (f)) as  
 611  $\sigma_0, a_1, \dots, a_k, \sigma_k \circ F_{1,x}(\epsilon) \circ a_{k+1}, \sigma_{k+1} \circ F_{1,y}(\epsilon) \circ S(\epsilon)$ , such that  $F_{1,x}(\epsilon) \stackrel{s_x}{\sim} F_{1,x}(\delta^{(k+1)})$ , where  
 612  $a_{k+1}, \sigma_{k+1}$  occurs after  $F_{1,x}(\delta^{(k+1)})$  and  $R(\epsilon)$  returns  $(x_1, y_1)$ .

613 Now, consider the finite execution  $\sigma_0, a_1, \dots, a_k, \sigma_k \circ F_{1,x}(\epsilon)$  at the end of which we append  
 614  $recv(m_y^{r_1})_{r_1, s_y}$  to create a finite execution of  $\mathcal{A}$  as  $\sigma_0, a_1, \dots, a_k, \sigma_k \circ F_{1,x}(\epsilon), recv(m_y^{r_1})_{r_1, s_y}$ .  
 615 Then, there exists an execution  $\epsilon'$  of  $\mathcal{A}$ , where the network delays the input actions at  
 616  $s_y$ . By the  $N$  and  $O$  properties,  $send(y)_{s_y, r_1}$  occurs in  $\epsilon'$ . Clearly, since  $F_{1,x}(\epsilon') \stackrel{s_x}{\sim} F_{1,x}(\epsilon)$ ,  
 617  $send(x)_{s_x, r_1}$  sends the same value in  $\epsilon$  and  $\epsilon'$ . Therefore,  $R_1$  returns  $(x_1, y)$ .

618 In  $\epsilon'$ , we denote the fragment that begins with  $recv(m_y^{r_1})_{r_1, s_y}$  and ends with  $send(y)_{s_y, r_1}$   
 619 by  $F_{1,y}(\epsilon')$  to have  $\epsilon'$  as  $\sigma_0, a_1, \dots, a_k, \sigma_k \circ F_{1,x}(\epsilon') \circ F_{1,y}(\epsilon')$ . Then, there exists an execution  
 620  $\delta^{(k)}$  of  $\mathcal{A}$ , which is an extension of  $\epsilon'$ . Clearly,  $\delta^{(k)}$  can be written as  $\sigma_0, a_1, \dots, a_k, \sigma_k \circ$   
 621  $F_{1,x}(\delta^{(k)}) \circ F_{1,y}(\delta^{(k)}) \circ S(\delta^{(k)})$ , where  $S(\delta^{(k)})$  is the tail of the extended execution. Clearly,  
 622  $F_{1,x}(\delta^{(k)}) \stackrel{s_x}{\sim} F_{1,x}(\epsilon')$ . Therefore,  $R_1(\delta^{(k)})$  returns  $x_1$  in  $\delta^{(k)}$ , which implies  $R_1(\delta^{(k)})$  must  
 623 return  $(x_1, y_1)$ . ◀

## 624 5.2 SNOW with C2C Communication

625 In this section, we state that SNOW is possible in the *multiple-writers single-reader* (MWSR)  
 626 setting when client-to-client communication is allowed. We consider a system that has  $\ell \geq 1$   
 627 writers with ids  $w_1, w_2, \dots, w_\ell \in \mathcal{W}$ , one reader  $r$ , and  $k \geq 1$  servers with ids  $s_1, s_2, \dots, s_k \in \mathcal{S}$ .  
 628 Client-to-client communication is allowed. The following theorem states that in such a  
 629 setting it is possible to have an algorithm that respects all SNOW properties. Due to space  
 630 constraints, we do not present the algorithm or the proof for the theorem, which can be  
 631 found in Appendix E.3.

► **Theorem 23.** *In the MWSR setting, there exists an algorithm  $A$  such that any well-formed and fair execution of  $A$  implements a wait-free transaction processing system,  $T$ , for objects of type  $\mathcal{O}_T$ , consisting of objects  $o_1, o_2, \dots, o_k$  at servers  $s_1, s_2, \dots, s_k$ , respectively; and  $T$  respects the SNOW properties and WRITE transactions are live.*

## 6 Conclusion

We revisited the SNOW Theorem and when it is possible for READ transactions to have the same latency as simple reads. We provided a new and more rigorous proof of the original result. We also closed several open questions that were either explicitly posed by the original work or that emerged from our increased rigor. We found that READ transactions can match the latency of simple reads when client-to-client communication is allowed and there is at most one reader. We found that they cannot and must have higher worst-case latency when client-to-client communication is disallowed or there are at least two readers. We also presented the first algorithms that provide bounded worst-case latency for read-only transactions in strictly serializable systems with WRITE transactions.

Finally, each of our algorithms assumes a set of readers and writers. Clients in real systems, however, dynamically switch between reading and writing data and thus are both. Does such a setting change the bounds when SNOW is possible?

## References

- 1 Hagit Attiya and Jennifer L. Welch. Sequential consistency versus linearizability. *ACM Trans. Comput. Syst.*, 12(2):91–122, 1994.
- 2 P.A. Bernstein, V. Hadzilacos, and N. Goodman. *Concurrency Control and Recovery in Database Systems*. Addison-Wesley Longman Publishing, 1987.
- 3 Eric A. Brewer. Towards robust distributed systems. In *Proc. Principles of Distributed Computing*, Jul 2000.
- 4 Nathan Bronson and et. al. TAO: Facebook’s distributed data store for the social graph. In *USENIX Annual Technical Conference (USENIX ATC 13)*, pages 49–60, 2013.
- 5 Phil Dixon. Shopzilla site redesign: We get what we measure. Velocity Conf. Talk, 2009.
- 6 James C. Corbett et. al. Spanner: Google’s globally-distributed database. In *Proc. OSDI*, Oct 2012.
- 7 Michael J. Fischer, Nancy A. Lynch, and Michael S. Paterson. Impossibility of distributed consensus with one faulty process. In *Proc. Prin. of Database Sys*, 1983.
- 8 Seth Gilbert and Nancy Lynch. Brewer’s conjecture and the feasibility of consistent, available, partition-tolerant web services. In *ACM SIGACT News*, Jun 2002.
- 9 J. Gray and A. Reuter. *Transaction Processing: Concepts and Techniques*. Morgan Kaufmann, San Mateo, CA, 1993.
- 10 M. P. Herlihy and Jeannette M. Wing. Linearizability: a correctness condition for concurrent objects. *ACM Transactions on Programming Languages and Systems (TOPLAS)*, 12(3):463–492, 1990.
- 11 Greg Linden. Make data useful. Stanford CS345 Talk, 2006.
- 12 Richard J. Lipton and Jonathan S. Sandberg. PRAM: A scalable shared memory. Technical Report TR-180-88, Princeton Univ., Dept. Comp. Sci., 1988.
- 13 Wyatt Lloyd, Michael J. Freedman, Michael Kaminsky, and David G. Anderson. Stronger Semantics for Low-Latency Geo-Replicated Storage. In *Proc. NSDI*, Apr 2013.
- 14 H. Lu, C. Hodsdon, K. Ngo, S Mu, and W. Lloyd. The SNOW theorem and latency-optimal read-only transactions. In *12th USENIX Symp. on Operating Sys. Design and Implementation (OSDI 16)*, pages 135–150, Savannah, GA, 2016.
- 15 N. A. Lynch. *Distributed Algorithms*. Morgan Kaufmann Publishers, 1996.

- 679 16 Christos H. Papadimitriou. The serializability of concurrent database updates. *Journal of*  
 680 *ACM*, 26(4):631–653, 1979.
- 681 17 E. Schurman and J. Brutlag. The user and business impact of server delays, additional bytes,  
 682 and HTTP chunking in web search. Velocity Conf. Talk, 2009.

## 683 A Algorithm Specification with I/O Automata

684 We model a distributed algorithm using the I/O Automata [15]. Here we limit our discussion  
 685 of I/O Automata to the relevant concepts, but for a detailed account the reader should  
 686 refer to [15]. An algorithm is a composition  $\mathcal{A}$  of a set of *automata* where each automaton  
 687  $A_i$  corresponds to a process (e.g., client or server) or a communication channel in the  
 688 system.  $A_i$  is defined in terms of a set of deterministic transition functions  $trans(A_i)$   
 689 (also called *actions*), which can be thought of as the algorithmic steps of  $A_i$ ; and a set  
 690 of states  $states(A_i)$ . An execution of  $\mathcal{A}$  is a sequence of alternating states and actions of  
 691  $A$ ,  $\sigma_0, a_1, \sigma_1, a_2, \sigma_2, \dots, \sigma_n$ . A state change, called a *step*, is a 3-tuple  $(\sigma_i, a_i, \sigma_{i+1})$ , with  
 692  $\sigma_i, \sigma_{i+1} \in states(A_i)$  and  $a_i \in trans(A_i)$ . The set of input actions is denoted by  $in(A_i)$ , e.g.,  
 693  $a_i \in in(A_i)$  is an input action if it receives a message. The set of output actions is denoted  
 694 by  $out(A_i)$ . The input and output actions are also called external actions, denoted by  $ext(A_i)$   
 695 and  $in(A_i) \cup out(A_i) = ext(A_i)$ . If an action  $a_i \notin ext(A_i)$ , then  $a_i$  is an internal action.  
 696 Communications between any two automata  $A_i$  and  $A_j$  is modeled by using channel automata  
 697  $Channel_{i,j}$  for sending messages from  $A_i$  to  $A_j$ ; and  $Channel_{j,i}$  for sending message from  
 698  $A_j$  to  $A_i$ . When  $A_i$  sends some message  $m$  to  $A_j$  the following sequence of actions occur:  
 699  $send(m)_{i,j}$  occurs at  $A_i$ , then  $send(m)_{i,j}$  followed by  $recv(m)_{i,j}$  occur at  $Channel_{i,j}$ ; then  
 700 finally,  $A_j$  receives  $m$  via the action  $recv_{i,j}(m)$ . In our model, the communication channels are  
 701 simple because we assume reliable communication between each pair of processes. Therefore,  
 702 we ignore the actions in the  $Channel_{i,j}$  and instead say  $send(m)_{i,j}$  occurs at  $A_i$  and then  
 703  $A_j$  receives  $m$  via the action  $recv_{i,j}(m)$ . An execution fragment  $\alpha$  can be either finite, i.e.,  
 704 having finite states, or infinite. If  $\alpha$  is a finite execution and  $\beta$  is an execution fragment, such  
 705 that  $\beta$  starts with the final state of  $\alpha$  then we use  $\alpha \circ \beta$  to denote the concatenation of  $\alpha$   
 706 and  $\beta$ . If  $\epsilon$  and  $\epsilon'$  are two execution fragments, such that they have the same sequence of  
 707 states at automaton  $A_i$ , i.e.,  $\epsilon|A_i = \epsilon'|A_i$ , then  $\epsilon$  and  $\epsilon'$  are indistinguishable at  $A_i$ , denoted  
 708 by  $\epsilon \stackrel{A_i}{\sim} \epsilon'$ . When the context is clear, we simply use  $\epsilon \sim \epsilon'$ .

709 For any execution of  $\mathcal{A}$ ,  $\sigma_0, a_1, \dots, a_k, \sigma_k \dots$ , where  $\sigma$ 's and  $a$ 's are states and actions, we  
 710 use the notation  $a_1, \dots, a_k \dots$  that shows only the actions while leaving out the states to  
 711 simplify notation.

## 712 B Some Useful I/OA results

713 Below we add some useful theorems are useful related to executions of a composed I/O  
 714 Automata [15].

715 ▶ **Theorem 24.** Let  $\{A_i\}_{i \in I}$  be a compatible collection of automata and let  $A = \Pi_{i \in I} A_i$ .  
 716 Suppose  $\alpha_i$  is an execution of  $A_i$  for every  $i \in I$ , and suppose  $\beta$  is a sequence of actions in  
 717  $ext(A)$  such that  $\beta|A_i = trace(\alpha_i)$  for every  $i \in I$ . Then there is an execution  $\alpha$  of  $A$  such  
 718 that  $\beta = trace(\alpha)$  and  $\alpha_i = \alpha|A_i$ , for every  $i \in I$ .

719 ▶ **Theorem 25.** Let  $A$  be any I/O automaton.

- 720 1. If  $\alpha$  is a finite execution of  $A$ , then there is a fair execution of  $A$  that starts with  $\alpha$ .  
 721 2. If  $\beta$  is a finite trace of  $A$ , then there is a fair trace of  $A$  that starts with  $\beta$ .

- 722 3. If  $\alpha$  is a finite execution of  $A$  and  $\beta$  is any finite or infinite sequence of input actions of  
 723  $A$ , then there is a fair execution  $\alpha \circ \alpha'$  of  $A$  such that the sequence of input actions in  $\alpha'$   
 724 is exactly  $\beta$ .  
 725 4. If  $\beta$  is a finite trace of  $A$  and  $\beta'$  is any finite or infinite sequence of input actions of  
 726  $A$ , then there is a fair execution  $\alpha \circ \alpha'$  of  $A$  such that  $\text{trace}(\alpha) = \beta$  and such that the  
 727 sequence of input actions in  $\alpha'$  is exactly  $\beta'$ .

728 The following useful claim is adopted from Chapter 16 of [15].

729  $\triangleright$  **Claim 26.** Suppose we have an automaton  $A = \Pi_{i=1}^k A_i$  where  $A$  is composed of the  
 730 compatible collection of automata  $A_i$ , where  $i \in \{1, \dots, k\}$ . Let  $\beta$  be a fair trace of  $A$  then  
 731 we define an irreflexive partial order  $\rightarrow_\beta$  on the actions of  $\beta$  as follows. If  $\pi$  and  $\phi$  are events  
 732 in  $\beta$ , with  $\pi$  preceding  $\phi$ , then we say  $\phi$  depends on  $\pi$ , which we denote as  $\pi \rightarrow_\beta \phi$ , if one of  
 733 the following holds:

- 734 1.  $\pi$  and  $\phi$  are actions at the same automaton;  
 735 2.  $\pi$  is some  $\text{send}(\cdot)_{j,i}$  at some  $A_j$  and  $\phi$  is some  $\text{recv}(\cdot)_{j,i}$  at  $A_i$ ; and  
 736 3.  $\pi$  and  $\phi$  are related by a chain of the relations of items 1. and 2.

737 Then if  $\gamma$  is a sequence obtained by reordering the events in  $\beta$  while preserving the  $\rightarrow_\beta$ , then  
 738  $\gamma$  is also a fair trace of  $A$ .

739  $\blacktriangleright$  **Theorem 27 ([15]).** Let  $\{A_i\}_{i \in I}$  be a compatible collection of automata and let  $A = \Pi_{i=1}^k A_i$ .  
 740 Suppose  $\alpha_i$  is a fair execution of  $A_i$  for every  $i \in I$ , and suppose  $\beta$  is a sequence of actions  
 741 in  $\text{ext}(A)$  such that  $\beta|A_i = \text{trace}(\alpha_i)$  for every  $i \in I$ . Then there is a fair execution  $\alpha$  of  $A$   
 742 such that  $\beta = \text{trace}(\alpha)$  and  $\alpha_i = \alpha|A_i$  for every  $i \in I$ .

## 743 **C** Technical Preliminaries

744  $\blacktriangleright$  **Lemma.** 3 Suppose  $\alpha$  is any execution of  $\mathcal{A}$  such that a  $\text{READ } R_i$  is in  $\alpha$ . Suppose the  
 745 execution fragment  $I_i(\alpha) \circ F_{i,x}(\alpha)^{(x_t)} \circ F_{i,y}(\alpha)^{(y_s)} \circ E(\alpha)^{(x_{t'}, y_{s'})}$  in  $\alpha$ , corresponds to  $R_i$ ,  
 746 where  $x_t, x_{t'} \in V_1$  and  $y_s, y_{s'} \in V_2$ , and  $s, s', t, t'$  are version identifiers then (i)  $s = s'$  and  
 747  $t = t'$  and (ii)  $s' = t'$ .

748 **Proof.** Suppose  $R_i$  is invoked at reader  $r_i$ . Then, via the action  $\text{send}(x_t)_{s_x, r_i}$ , in execution  
 749 fragment  $F_{i,x}(\alpha)^{(x_t)}$ , server  $s_1$  sends the value  $x_t$  to  $r_i$ , which is received at  $r_i$  through the  
 750 action  $\text{recv}(x_{t'})_{s_1, r_i}$  in  $E_i(\alpha)^{(x_{t'}, y_{s'})}$ . By the assumptions of the reliable channel automata in  
 751 our model, we have  $x_t = x_{t'}$ , i.e.,  $t = t'$ . Similar argument for  $F_{i,y}(\alpha)^{(y_s)}$  and  $E_i(\alpha)^{(x_{t'}, y_{s'})}$   
 752 leads us to conclude  $s = s'$ . Next,  $R_i$  responds with  $(x_{t'}, y_{s'})$ , which implies by the S property  
 753 for executions of  $\mathcal{A}$  that  $x_{t'}$  and  $y_{s'}$  must correspond to the same version, i.e.,  $s' = t'$ .  $\blacktriangleleft$

754 Note that the above results hold even if there are any other execution fragments, that do  
 755 not contain any actions at  $r_i$ ,  $s_1$  or  $s_y$ , in-between the  $I_i$ ,  $F_{i,j}$  and  $E_i$  execution fragments.

756  $\blacktriangleright$  **Corollary 28.** Suppose  $\alpha$  is any execution of  $\mathcal{A}$  such that a  $\text{READ } R_i$  is in  $\alpha$ . Suppose the  
 757 execution fragment  $I_i(\alpha) \circ X_1 \circ F_{i,x}(\alpha)^{(x_t)} \circ X_2 \circ F_{i,y}(\alpha)^{(y_s)} \circ X_3 \circ E_i(\alpha)^{(x_{t'}, y_{s'})}$  in  $\alpha$ , corresponds  
 758 to  $R_i$ , where  $x_t, x_{t'} \in V_1$  and  $y_s, y_{s'} \in V_2$ ,  $X_1, X_2, X_3$  are some execution fragments that do  
 759 not contain any action at  $r_i$ ,  $s_1$  or  $s_y$ , and  $s, s', t, t'$  are version identifiers then (i)  $s = s'$   
 760 and  $t = t'$  and (ii)  $s' = t'$ .

761 **Proof.** The constraint (i) in the statement can be derived from the fact that  $\mathcal{A}$  satisfies the  
 762 O property, which implies that any version returned by  $R_i$  for object value  $o_1$  (or  $o_2$ ) and  
 763 this must be the only version sent by the server  $s_1$  (or  $s_y$ ) to  $r$ . Constraint (ii) from the S  
 764 property of  $\mathcal{A}$ .  $\blacktriangleleft$

765 ► **Lemma. 4** (Commuting fragments) Let  $\alpha$  be an execution of  $\mathcal{A}$ . Suppose  $G_1(\alpha)$  and  $G_2(\alpha)$   
 766 are any execution fragments in  $\alpha$  such that all actions in each fragment occur only at one  
 767 automaton and either (a) none of the fragments contain input actions, or (b) at least one  
 768 of the fragments have no external actions. Suppose  $G_1(\alpha)$  and  $G_2(\alpha)$  occur at two distinct  
 769 automata and the execution fragment  $G_1(\alpha) \circ G_2(\alpha)$  occurs in  $\alpha$ . Then there exists an  
 770 execution  $\alpha'$  of  $\mathcal{A}$ , where the execution fragment  $G_2(\alpha) \circ G_1(\alpha)$  appears in  $\alpha'$ , such that (i)  
 771  $G_1(\alpha) \sim G_1(\alpha')$  and  $G_2(\alpha) \sim G_2(\alpha')$  (ii) the prefix in  $\alpha$  before  $G_1(\alpha) \circ G_2(\alpha)$  is identical  
 772 to the prefix in  $\alpha'$  before  $G_1(\alpha') \circ G_2(\alpha')$ ; and (ii) the suffix in  $\alpha$  after  $G_1(\alpha) \circ G_2(\alpha)$  is  
 773 identical to the suffix in  $\alpha'$  after the execution fragment  $G_2(\alpha') \circ G_1(\alpha')$ .

774 **Proof.** This is clear because the network can move the actions in  $G_2$  to occur before  $G_1$   
 775 at their respective automata, and because either (a) none of the fragments have any input  
 776 action or (b) at least one of them has no external actions, and hence the actions in one of  
 777 these fragments cannot affect the actions in the other fragment. ◀

778 ► **Lemma. 5** (Indistinguishability) Let  $\alpha$  and  $\beta$  be executions of  $\mathcal{A}$  and let  $R$  be any READ  
 779 transaction. Then (i) if  $F_x(\alpha) \stackrel{s_x}{\sim} F_x(\beta)$  then both  $R(\alpha)$  and  $R(\beta)$  respond with the same  
 780 value  $x$  at  $s_x$ ; and (ii) if  $F_y(\alpha) \stackrel{s_y}{\sim} F_y(\beta)$  then both  $R(\alpha)$  and  $R(\beta)$  respond with the same  
 781 value  $y$  at  $s_y$ .

782 **Proof.** Suppose  $R$  is invoked at some reader  $r$ . Let  $j \in \{1, 2\}$  and suppose the fragments  
 783  $F_j(\alpha)$  and  $F_j(\beta)$  appears in  $\alpha$  and  $\beta$  respectively, where in  $F_j(\alpha)$  server  $s_j$  sends  $v_j \in V_j$  to  $r$ .  
 784 Then  $R(\alpha)$  must return  $v_j$  for object  $o_j$  by the O property of  $\mathcal{A}$ . Then since  $F_j(\alpha) \stackrel{s_j}{\sim} F_j(\beta)$   
 785 then in  $F_j(\beta)$  the server  $s_j$  must also send  $v_j$  to  $r$ , therefore, both  $R(\alpha)$  and  $R(\beta)$  must  
 786 return value  $v_j$  for  $o_j$ . ◀

787 ► **Lemma. 6** If any finite execution of  $\mathcal{A}$  ends with  $INV(R)$ , for a READ transaction  
 788  $R$  then there exists an extension  $\alpha$  which is a fair execution of  $\mathcal{A}$  and is of the form  
 789  $P(\alpha) \circ I(\alpha) \circ F_1(\alpha)^{(x)} \circ F_2(\alpha)^{(y)} \circ E(\alpha)^{(x,y)} \circ S(\alpha)$ , where  $P(\alpha)$  is the prefix and  $S(\alpha)$  denotes  
 790 the rest of the execution.

791 **Proof.** Consider a finite execution of  $\mathcal{A}$  that end with  $INV(R)$ , which occurs at some  
 792 reader  $r$ , then the network induces the execution fragment  $I(\alpha)$  by delaying all actions,  
 793 except the internal and output actions at  $r$ , between the actions  $INV(R)$  and the later  
 794 of the actions  $send(m_x^r)_{r,s_x}$  and  $send(m_y^r)_{r,s_y}$ . Next, the network delivers  $m_x^r$  at  $s_x$  (via  
 795 the action  $recv(m_x^r)_{r,s_x}$ ) and delays all actions, other than internal and output actions  
 796 at  $s_x$ , until  $s_x$  responds with  $x$ , via  $send(x)_{s_x,r}$ ; we identify this execution fragment as  
 797  $F_1(\alpha)^{(x)}$ . Subsequently, in a similar manner, the network delivers the message  $m_y^r$  and delays  
 798 appropriate actions to induce the execution fragment  $F_2(\alpha)^{(y)}$ . Finally, the network delivers  
 799 the values  $x$  and  $y$  to  $r$  (via the events  $recv(x)_{s_x,r}$  and  $recv(y)_{s_y,r}$ ), and delays all actions at  
 800 other automata until  $R$  completes with action  $RESP(R)$  by returning  $(x, y)$ . As a result, we  
 801 arrive at a fair execution of  $\mathcal{A}$  of the form  $I(\alpha) \circ F_1(\alpha)^{(x)} \circ F_2(\alpha)^{(y)} \circ E(\alpha)^{(x,y)} \circ S(\alpha)$ . ◀

## 802 **D** Proof of the SNOW Theorem 3-clients

803 ► **Lemma. 9** [Existence of  $\alpha_2$ ] There exists fair execution  $\alpha_2$  of  $\mathcal{A}$  that contains transactions  
 804  $W$ ,  $R_1$  and  $R_2$  and can be written in the form  $P_{k+1} \circ R_1^{(x_1,y_1)} \circ R_2^{(x_1,y_1)} \circ S$ , where both  $R_1$   
 805 and  $R_2$  return  $(x_1, y_1)$ .

806 **Proof.** We can construct a fair execution  $\alpha_2$  of  $\mathcal{A}$  as follows. Consider the prefix  $a_1, \dots, a_{k+1} \circ$   
 807  $R_1(\alpha_1)^{(x_1,y_1)}$  of the execution  $\alpha_1$ , from Lemma 8. At the end of this prefix, the network

invokes  $R_2$ . Now, by Lemma 6, due to  $INV(R_2)$  there is an extension of the prefix of the form  $a_1, \dots, a_{k+1} \circ R_1(\alpha_1)^{(x_1, y_1)} \circ I(\alpha) \circ F_{1,x}(\alpha)^{(x)} \circ F_{1,y}(\alpha)^{(y)} \circ E(\alpha)^{(x,y)}$ . By the S property, we have  $x = x_1$  and  $y = y_1$ . Therefore,  $\alpha_2$  (Fig. 1b) can be written in the form  $P_{k+1} \circ R_1^{(x_1, y_1)} \circ R_2^{(x_1, y_1)} \circ S$ , where  $S$  is the rest of the execution.  $\blacktriangleleft$

In the following lemma, we show that we can create a fair execution  $\alpha_4$ , of  $\mathcal{A}$ , where  $F_{2,y}$  occurs immediately before  $E_1^{(x_1, y_1)}$ , while  $R_1$  and  $R_2$  both return  $(x_1, y_1)$ .

► **Lemma.** 11 [Existence of  $\alpha_4$ ] *There exists fair execution  $\alpha_4$  of  $\mathcal{A}$  that contains transactions  $W$ ,  $R_1$  and  $R_2$  and can be written in the form  $P_k \circ I_2 \circ a_{k+1} \circ I_1 \circ F_{1,x} \circ F_{1,y} \circ F_{2,y} \circ E_1 \circ F_{2,x} \circ E_2 \circ S$ , where both  $R_1$  and  $R_2$  return  $(x_1, y_1)$ .*

**Proof.** We start with an execution  $\alpha_3$ , as in Lemma 10, and apply Lemma 4 twice.

First, by Lemma 4, we know there exists a fair execution  $\alpha'$  of  $\mathcal{A}$  where  $F_{2,x}$  (identify as  $G_1$ ) and  $F_{2,y}$  (identify as  $G_2$ ) are interchanged since actions of  $F_{2,x}$  occurs solely at  $s_1$  and those of  $F_{2,y}$  at  $s_y$ , and  $F_{2,x}$  and  $F_{2,y}$  return  $x_1$  and  $y_1$ , respectively, to  $r_2$ .

Next, by Lemma 4 there is fair execution of  $\mathcal{A}$ , say  $\alpha_4$  where the fragments  $E_1$  (identify as  $G_1$ ) and  $F_{2,y}$  (identify as  $G_2$ ) are interchanged, with respect to  $\alpha'$ , because the actions in  $E_1$  occur at  $r_1$  and those of  $F_{2,y}$  at  $s_y$ . Furthermore,  $\alpha_4$  can be written in the form  $P_k \circ I_2 \circ a_{k+1} \circ I_1 \circ F_{1,x}^{(x_1)} \circ F_{1,y}^{(y_1)} \circ F_{2,y}^{(y_1)} \circ E_1^{(x_1, y_1)} \circ F_{2,x}^{(x_1)} \circ E_2^{(x_1, y_1)} \circ S$ .  $\blacktriangleleft$

In the following lemma, starting from  $\alpha_6$  in Lemma 13 we create a fair execution  $\alpha_7$  for  $\mathcal{A}$  where  $F_{2,x}$  appears before  $F_{1,y} \circ E_1$ , where  $R_1$  returns  $(x_0, y_0)$  and  $R_2$  returns  $(x_1, y_1)$ . At high level, we will be working on moving the execution fragments of  $R_2$  forward, a little at a time, until finally we have  $R_2$  finishing before  $R_1$  starts. This simply uses commutativity since the actions in the swapped execution fragments occurs at different automata.

► **Lemma.** 14 [Existence of  $\alpha_7$ ] *There exists fair execution  $\alpha_7$  of  $\mathcal{A}$  that contains transactions  $W$ ,  $R_1$  and  $R_2$ , and can be written in the form  $P_k \circ I_2 \circ I_1 \circ F_{1,x} \circ F_{2,y} \circ F_{2,x} \circ F_{1,y} \circ E_1 \circ E_2 \circ S$  where  $R_1$  returns  $(x_0, y_0)$  and  $R_2$  returns  $(x_1, y_1)$ .*

**Proof.** This result is proved by applying the result of Lemma 4 to the fair execution created in Lemma 13. Suppose,  $\alpha_6$  (Fig. 1b) is a fair execution as in Lemma 13, where in the execution fragment  $E_1^{(x_0, y_0)} \circ F_{2,x}$  we identify  $E_1^{(x_0, y_0)}$  as  $G_1$  and  $F_{2,x}^{(y_0)}$  as  $G_2$ . The actions of  $G_1$  and  $G_2$  occur at two distinct automata, therefore, we can use the result of Lemma 4, to argue that there exists a fair execution  $\alpha'$  of  $\mathcal{A}$  that contains the execution fragment  $F_{2,x} \circ E_1^{(x_0, y_0)}$ , and  $\alpha_6$  and  $\alpha'$  are identical in the prefixes and suffixes corresponding to  $G_1$  and  $G_2$ .

Now,  $\alpha'$  contains  $F_{1,y} \circ F_{2,x}$ , where the actions in  $F_{1,y}$  (identified as  $G_1$ ) and  $F_{2,x}$  (identify as  $G_2$ ) occur at distinct automata. Hence, by Lemma 4 there exists an execution  $\alpha_7$  of the form  $P_k \circ I_2 \circ I_1 \circ F_{1,x}^{(x_0)} \circ F_{2,y}^{(y_1)} \circ F_{2,x}^{(x_1)} \circ F_{1,y}^{(y_0)} \circ E_1^{(x_0, y_0)} \circ E_2^{(x_1, y_1)} \circ S$ .  $\blacktriangleleft$

► **Lemma.** 15 [Existence of  $\alpha_8$ ] *There exists fair execution  $\alpha_8$  of  $\mathcal{A}$  that contains transactions  $W$ ,  $R_1$  and  $R_2$  and can be written in the form  $P_k \circ I_2 \circ F_{2,y} \circ I_1 \circ F_{1,x} \circ F_{2,x} \circ F_{1,y} \circ E_1 \circ E_2 \circ S$ , where  $R_1$  returns  $(x_0, y_0)$  and  $R_2$  returns  $(x_1, y_1)$ .*

**Proof.** Consider the fair execution  $\alpha_7$  of  $\mathcal{A}$  as in Lemma 14. In the context of of Lemma 4, in  $\alpha_7$  (Fig. 1b) the actions in  $F_{1,x}$  (identify as  $G_1$ ) occur at  $s_1$  and those in  $F_{2,y}$  (identify as  $G_2$ ) at  $s_y$ . Then by Lemma 4 there exists a fair execution  $\alpha'$  of  $\mathcal{A}$ , of the form  $P_k \circ I_2 \circ I_1 \circ F_{2,y} \circ F_{1,x} \circ F_{2,x} \circ F_{1,y} \circ E_1 \circ E_2 \circ S$ , where  $F_{2,y}$  and  $F_{1,x}$  are interchanged.



Since actions in  $F_{2,y}$  (identify as  $G_1$ ) occur at  $s_y$  and those in  $I_1$  (identify as  $G_1$ ) occur at  $r_1$  then by Lemma 4 there is a fair execution of  $\mathcal{A}$ ,  $\alpha_8$  where  $F_{2,y}$  appear before  $I_1$ , i.e., of the form  $P_k \circ I_2 \circ F_{2,y} \circ I_1 \circ F_{1,x} \circ F_{2,x} \circ F_{1,y} \circ E_1 \circ E_2 \circ S$ , where  $F_{2,y}$  and  $I_1$  are interchanged.

By (ii) of Lemma 5 we have  $F_{2,x}(\alpha') \stackrel{s_1}{\sim} F_{2,x}(\alpha_8)$  hence  $F_{2,x}$  sends  $x_1$  and  $F_{1,x}$  and  $F_{1,y}$  sends  $x_0$  and  $y_0$ , respectively. So considering these returned values we have  $\alpha_8$  (Fig. 1b) in the form as stated in the lemma.  $\blacktriangleleft$

► **Lemma 29.** 16 [Existence of  $\alpha_9$ ] There exists fair execution  $\alpha_9$  of  $\mathcal{A}$  that contains transactions  $W$ ,  $R_1$  and  $R_2$  and can be written in the form  $P_k \circ I_2 \circ F_{2,y}^{(y_1)} \circ I_1 \circ F_{2,x}^{(x_1)} \circ F_{1,x}^{(x_0)} \circ F_{1,y}^{(y_0)} \circ E_1^{(x_0,y_0)} \circ E_2^{(x_1,y_1)} \circ S$  where  $R_1$  returns  $(x_0, y_0)$  and  $R_2$  returns  $(x_1, y_1)$ .

**Proof.** In  $\alpha_8$  from Lemma 15, all the actions in  $I_1$  occur at  $r_1$ ; those in  $F_{1,x}$  occur at  $s_1$ ; and the actions in  $F_{2,x}$  occur only at  $s_1$ . Note that actions of both execution fragments  $F_{2,x}$  and  $F_{1,x}$  occur at  $r_1$ . Consider the prefix of  $\alpha_8$  that ends with  $I_1$  then suppose the network extends this prefix by adding an execution fragment of the form  $F_{2,x} \circ F_{1,x}$  as follows. First note that the actions  $send(m_x^{r_2})_{r_2,s_1}$  and  $send(m_x^{r_1})_{r_1,s_1}$  appears in the prefix but do not have corresponding *recv* actions. The network places action  $recv(m_x^{r_2})_{r_2,s_1}$ , and allows an execution fragment of the form  $F_{2,x}$  to appear. Now, immediately after this the network further extends it with an execution fragment of the form  $F_{1,x}$  by placing action  $recv(m_x^{r_1})_{r_1,s_1}$ . Next the fragment  $F_{1,y}$  is added and is the same as  $F_{1,y}(\alpha_8)$ . This last step can be argued by the fact that none of the actions in  $F_{1,y}$  can be affected by any of the output actions at  $F_{2,x}$  and  $F_{1,x}$ . Note that a careful argument can be done by using Theorem 27 to conclude the same. Following this the network allows the rest of the execution by adding an execution fragment of the form  $E_1 \circ E_2 \circ S$ . The resulting fair execution is of the form  $P_k \circ I_2 \circ F_{2,y}^{(y_1)} \circ I_1 \circ F_{2,x} \circ F_{1,x} \circ F_{1,y}^{(y_0)} \circ E_1 \circ E_2 \circ S$ , where we retained the values wherever it is known, and we denote this execution by  $\alpha_9$ .

Now, we argue about the return values in  $\alpha_9$ . Applying Lemma 5 to  $R_2$  and  $F_{2,y}$  implies that  $R_2$  returns  $(x_1, y_1)$ . Similarly, applying Lemma 5 to  $R_1$  and  $F_{1,y}$  implies that  $R_1$  must return  $(x_0, y_0)$  in  $\alpha_9$ .  $\blacktriangleleft$

## E Two Client Open Question

### E.1 No SNOW Without C2C Messages

► **Lemma.** 19 There exists a finite execution  $\alpha$  of  $\mathcal{A}$  that contains transactions  $R_1(\alpha)$  and  $W(\alpha)$  where  $INV(R_1)$  appears after  $RESP(W)$  and the following conditions hold:

- (i) The actions  $send(m_x^{r_1})_{r_1,s_1}$  and  $send(m_y^{r_1})_{r_1,s_y}$  appear consecutively in  $trace(\alpha)|_{r_1}$ ; and
- (ii)  $\alpha$  contains the execution fragment  $F_{1,x}(\alpha)$ .

**Proof.** Consider a finite execution fragment of  $\mathcal{A}$  with a completed transaction  $W$ , where after  $W$  completes the network invokes  $R_1$ , i.e.,  $INV(R_1)$  occurs. Note that each of the read operations  $op_1^{r_1}$  and  $op_2^{r_1}$ , in  $R_1$ , can be invoked by the network at any point in the execution. Following  $INV(R_1)$ , the network introduces the invocation action  $inv(op_1^{r_1})$ ; by the O property of the read operations of  $\mathcal{A}$  the action  $send(m_x^{r_1})_{r_1,s_1}$  eventually occurs. Next, the network introduces  $inv(op_2^{r_1})$  and also, delays the arrival of  $m_x^{r_1}$  until action  $send(m_y^{r_1})_{r_1,s_y}$  eventually occurs, which must occur in accordance with the property O of read operations. Let us call this finite execution  $\alpha^0$ .

Next, suppose at the end of  $\alpha^0$  the network delivers the message  $m_x^{r_1}$ , which has been delayed so far, via the action  $recv(m_x^{r_1})_{r_1,s_1}$ , at  $s_1$ , but it delays any other input actions at

893  $s_1$ . Note that by the  $N$  property of read operations  $s_1$  eventually responds with  $send(x)_{s_1, r_1}$ ,  
 894 with one value  $x$  by  $O$  property, where  $x = x_1$  by the  $S$  property, since  $R_1$  begins after  $W$   
 895 completes. Let us call this execution  $\alpha$ . Note that  $\alpha$  satisfies conditions (i) and (ii) by the  
 896 design of the execution.  $\blacktriangleleft$

897 **► Lemma. 20** *There exists an execution  $\beta$  of  $\mathcal{A}$  that contains transactions  $R_1$  and  $W$  where*  
 898  *$INV(R_1)$  appears after  $RESP(W)$  and the following conditions hold:*

- 899 (i) *The actions  $send(m_x^{r_1})_{r_1, s_1}$  and  $send(m_y^{r_1})_{r_1, s_y}$  appear consecutively in  $trace(\beta)|r_1$ ; and*  
 900 (ii)  *$\beta$  contains the execution fragment  $F_{1,x}(\beta) \circ F_{1,y}(\beta)$ .*

901 **Proof.** Consider the execution  $\alpha$  of  $\mathcal{A}$  as constructed in Lemma 19. At the end of the  
 902 execution fragment  $\alpha$ , the network delivers the previously delayed message  $m_y^{r_1}$ , which is  
 903 sent via the action  $send(m_y^{r_1})_{r_1, s_y}$ , by introducing the action  $recv(m_y^{r_1})_{r_1, s_y}$ . The network  
 904 then delays any other input action in  $\mathcal{A}$ . By the  $N$  property, server  $s_y$  must respond to  $r_1$ ,  
 905 with some value  $y$ , and hence the output action  $send(y)_{s_y, r_1}$  must eventually occur at  $s_y$ .  
 906 Let us call this finite execution as  $\beta$ . Note that  $\beta$  satisfies the properties (i) and (ii) in the  
 907 statement of the lemma.  $\blacktriangleleft$

908 **► Lemma. 21** *There exists a fair execution  $\gamma$  of  $\mathcal{A}$  with transactions  $R_1$  and  $W$  where the*  
 909 *action  $INV(R_1)$  appears before  $INV(W)$  and  $RESP(R_1)$  appears after  $RESP(W)$ , and the*  
 910 *following conditions hold for  $\gamma$ :*

- 911 (i) *The actions  $send(m_x^{r_1})_{r_1, s_x}$  and  $send(m_y^{r_1})_{r_1, s_y}$  appear before  $INV(W)$  and they appear*  
 912 *consecutively in  $trace(\gamma)|r_1$ ;*  
 913 (ii)  *$\gamma$  contains the execution fragment  $F_{1,x}(\gamma) \circ F_{1,y}(\gamma)$ ; and*  
 914 (iii) *action  $RESP(W)$  occurs before  $F_{1,x}(\gamma)$ .*

915 **Proof.** Consider the execution  $\beta$  of  $\mathcal{A}$  as in Lemma 20. Note that  $\beta$  is an execution of the  
 916 composed automaton  $\mathcal{A} (\equiv S_x \times r_1)$ . In  $\beta$ , the actions  $send(m_x^{r_1})_{r_1, s_x}$  and  $send(m_y^{r_1})_{r_1, s_y}$   
 917 occur at  $r_1$ ; and following that, the actions  $recv(m_x^{r_1})_{r_1, s_x}$ ,  $recv(m_y^{r_1})_{r_1, s_y}$ ,  $send(x)_{s_x, r_1}$  and  
 918  $send(y)_{s_y, r_1}$  occur at  $S_x$ . Consider the executions  $\alpha_{r_1} \equiv \beta|_{r_1}$  and  $\alpha_{S_x} \equiv \beta|_{S_x}$ . Let  $s_\beta$  denote  
 919  $trace(\beta)$ .

920 In  $s_\beta$ ,  $send(m_x^{r_1})_{r_1, s_x}$ ,  $send(m_y^{r_1})_{r_1, s_y}$  appear after  $RESP(W)$ , as in  $trace(\beta)$ . Let  $s'_\beta$  be  
 921 the sequence of external actions of  $S_y$  which we construct from  $s_\beta$  by moving  $send(m_x^{r_1})_{r_1, s_x}$ ,  
 922  $send(m_y^{r_1})_{r_1, s_y}$  before  $INV(W)$ , which is also an external action of  $\mathcal{A}$ , and leaving the rest of  
 923 the actions in  $s_\beta$  as it is.

924 In  $\beta$ ,  $INV(R_1)$ ,  $recv(x)_{s_x, r_1}$  and  $recv(y)_{s_y, r_1}$  are the only input actions at  $r_1$ , therefore,  
 925  $s'_\beta|_{r_1} = trace(\alpha_{r_1})$ . On the other hand,  $recv(m_x^{r_1})_{r_1, s_x}$ ,  $recv(m_y^{r_1})_{r_1, s_y}$  are the only input  
 926 actions at  $s_x$ , therefore,  $s'_\beta|_{S_1} = trace(\alpha_{S_1})$ . Now, by Theorem 24, there exists an  
 927 execution  $\gamma$  of  $\mathcal{A}$  such that,  $s'_\beta = trace(\gamma)$  and  $\alpha_{r_1} = \gamma|_{r_1}$  and  $\alpha_{S_x} = \gamma|_{S_x}$ . Therefore, in  $\gamma$ ,  
 928  $send(m_x^{r_1})_{r_1, s_x}$ ,  $send(m_y^{r_1})_{r_1, s_y}$  appear before  $INV(W)$  (condition (i)) and since  $s'_\beta = trace(\gamma)$   
 929 condition (ii) holds. Conditions (iii) holds trivially.  $\blacktriangleleft$

930 **► Lemma. 22** *There exists a fair execution  $\eta$  of  $\mathcal{A}$  that contains transactions  $R_1$  and  $W$*   
 931 *where  $INV(R)$  appears before  $INV(W)$ ;  $RESP(R_1)$  appears after  $RESP(W)$  and the following*  
 932 *conditions hold for  $\eta$ :*

- 933 (i)  *$\eta$  can be written in the form  $P(\eta) \circ F_{1,x}(\eta) \circ F_{1,y}(\eta) \circ S(\eta)$ , for some  $P(\eta)$  and  $S(\eta)$ ;*  
 934 (ii) *The actions  $send(m_x^{r_1})_{r_1, s_1}$  and  $send(m_y^{r_1})_{r_1, s_y}$  appear before  $INV(W)$  and they appear*  
 935 *consecutively in  $trace(\eta)|r_1$ ;*  
 936 (iii) *action  $RESP(W)$  occurs before  $F_{1,x}(\eta)$ ; and*

937 (iv)  $R_1(\eta)$  returns  $(x_1, y_1)$ .

938 **Proof.** Let  $\gamma$  be a fair execution of  $\mathcal{A}$ , as described in Lemma 21. Let  $\gamma^0$  be the execution  
 939 fragment of  $\gamma$  up to the action  $send(y)_{s_y, r_1}$ . Now, by Theorem 25 (1), there exists a fair  
 940 execution  $\gamma^0 \circ \mu$ , of  $\mathcal{A}$ , where  $\mu$  denotes the extended portion of the execution.

941 Clearly, by the N and O properties, the actions  $resp(op_1^{r_1})$  and  $resp(op_2^{r_1})$  must eventually  
 942 occur in  $\gamma^0 \circ \mu$ . Now, identify  $\eta$  as  $\gamma^0 \circ \mu$ , where  $P(\eta) \circ F_{1,x}(\eta) \circ F_{1,y}(\eta)$  is  $\gamma^0$ , and  $\mu$  is  $S(\eta)$ ,  
 943 thereby, proving condition (i).

944 Note the condition (ii) is satisfied by  $\eta$  because  $RESP(W)$  appears in  $P(\eta)$ , therefore,  
 945 the fair execution  $\gamma$  is equivalent to the execution fragment of  $P(\eta)$  up to the event  $INV(W)$ ,  
 946 and also,  $\gamma$  satisfies condition (ii) as stated in Lemma 21.

947 Condition (iii) is true because  $F_{1,x}(\eta)$  begins with action  $recv(m_x^{r_1})_{r_1, s_1}$ , which occurs  
 948 after  $RESP(W)$ . Condition (iv) is satisfied by  $\eta$  because  $\eta$  is an extension of  $\gamma$  and due to  
 949 the result of Lemma ??.

## 950 E.2 Condition for Proving Strict Serializability

951 In this section, we derive a useful property for executions of algorithms that implement  
 952 objects of data type  $\mathcal{O}_T$  that will later help us show the *strict serializability* property (S  
 953 property) of algorithms presented in later sections.

954 Although the strict serializability property in transaction-processing systems is a well-  
 955 studied topic, the specific setting considered in this paper is much simpler. Therefore, this  
 956 allows us to derive simpler conditions to prove the safety of these algorithms. A wide range  
 957 of transaction types and transaction processing systems are considered in the literature.  
 958 For example, in [16], Papadimitriou defined the strict serializability conditions as a part  
 959 of developing a theory for analyzing transaction processing systems. In this work, each  
 960 transaction  $T$  consists of a set of write operations  $W$ , at individual objects, and a set of  
 961 read operations  $R$  from individual objects, where the operations in  $W$  must complete before  
 962 the operations in  $R$  execute. Other types of transaction processing systems allow nested  
 963 transactions [9], where the transactions may contain sub-transactions [2] which may further  
 964 contain a mix of read or write operations, or even child-transactions. In most transaction  
 965 processing systems considered in the literature, transactions can be *aborted* so as to handle  
 966 failed transactions. As a result, the serializability theories are developed while considering  
 967 the presence of aborts. However, in our system, we do not consider any abort, nor any client  
 968 or server failures. A transaction in our system is either a set of independent writes or a set  
 969 of reads with all the reads or writes in a transaction operating on different objects. Such  
 970 simplifications allow us to formulate an equivalent condition for the execution of an algorithm  
 971 to prove the S property of such algorithms while implementing an object of data type  $\mathcal{O}_T$ .

972 We note that an execution of a variable of type  $\mathcal{O}_T$  is a finite se-  
 973 quence  $\mathbf{v}_0, INV_1, RESP_1, \mathbf{v}_1, INV_2, RESP_2, \mathbf{v}_2, \dots, \mathbf{v}_r$  or an infinite sequence  
 974  $\mathbf{v}_0, INV_1, RESP_1, \mathbf{v}_1, INV_2, RESP_2, \mathbf{v}_2, \dots$ , where  $INV$ 's and  $RESP$ 's are invocations and  
 975 responses, respectively. The  $\mathbf{v}_i$ 's are tuples of the the form  $(v_1, v_2, \dots, v_k) \in \prod_{i=1}^q V_i$ , that  
 976 corresponds to the latest values stored across the objects  $o_1, o_2 \dots o_k$ , and the values in  $\mathbf{v}_0$   
 977 are the initial values of the objects. Any adjacent quadruple such as  $\mathbf{v}_i, INV_{i+1}, RESP_{i+1},$   
 978  $\mathbf{v}_{i+1}$  is consistent with the  $f$  function for an object of type  $\mathcal{O}_T$  (see Section ??) . Now, the  
 979 safety property of such an object is a trace that describes the correct response to a sequence  
 980 of  $INV$ 's when all the transactions are executed sequentially. The *strict serializability* of  $\mathcal{O}_T$   
 981 says that each trace produced by an execution of  $\mathcal{O}_T$  with concurrent transactions appears  
 982 as some trace of  $\mathcal{O}_T$ . We describe this below in more detail.

► **Definition 30** (Strict-serializability). *Let us consider an execution  $\beta$  of an object of type  $\mathcal{O}_T$ , such that the invocations of any transaction at any client respects the well-formedness property. Let  $\Pi$  denote the set of complete transactions in  $\beta$  then we say  $\beta$  satisfies the strict-serializability property for  $\mathcal{O}_T$  if the following are possible:*

- (i) *For every complete READ or WRITE transaction  $\pi$  we insert a point (serialization point)  $\pi_*$  between the actions  $INV(\pi)$  and  $RESP(\pi)$ .*
- (ii) *We select a set  $\Phi$  of incomplete transactions in  $\beta$  such that for each  $\pi \in \Phi$  we select a response  $RESP(\pi)$ .*
- (iii) *For each  $\pi \in \Phi$  we insert  $\pi_*$  somewhere after  $INV(\pi)$  in  $\beta$ , and remove the  $INV$  for the rest of the incomplete transactions in  $\beta$ .*
- (iv) *If we assume for each  $\pi \in \Pi \cup \Phi$  both  $INV(\pi)$  and  $RESP(\pi)$  to occur consecutively at  $\pi_*$ , with the interval of the transaction shrunk to  $\pi_*$ , then the sequence of transactions in this new trace is a trace of an object of data type  $\mathcal{O}_T$ .*

Now, we consider any automaton  $\mathcal{B}$  that implements an object of type  $\mathcal{O}_T$ , and prove a result that serves us an equivalent condition for proving the strict serializability property of  $\mathcal{B}$ . Any trace property  $P$  of an automaton is a *safety property* if the set of executions in  $P$  is non-empty; *prefix-closed*, meaning any prefix of an execution in  $P$  is also in  $P$ ; and *limit-closed*, i.e., if  $\beta_1, \beta_2, \dots$  is any infinite sequence of executions in  $P$  is such that  $\beta_i$  is prefix of  $\beta_{i+1}$  for any  $i$ , then the limit  $\beta$  of the sequence of executions  $\{\beta_i\}_{i=0}^\infty$  is also in  $P$ . From Theorem 13.1 in [15], we know that the trace property, which we denote by  $P_{SC}$ , of any well-formed execution of  $\mathcal{B}$  that satisfies the strict-serializability property is a safety property. Moreover, from Lemma 13.10 in [15] we can deduce that if every execution of  $\mathcal{B}$  that is well-formed and failure-free, and also contains no incomplete transactions, satisfies  $P_{SC}$ , then any well-formed execution of  $\mathcal{B}$  that can possibly have incomplete transactions is also in  $P_{SC}$ . Therefore, in the following lemma, which gives us an equivalent condition for the strict serializability property of an execution  $\beta$ , we consider only executions without any incomplete transactions. The lemma is proved in a manner similar to Lemma 13.16 in [15], for atomicity guarantee of a single multi-reader multi-writer object.

► **Lemma 31.** *Let  $\beta$  be an execution (finite or infinite) of an automaton  $\mathcal{B}$  that implements an object of type  $\mathcal{O}_T$ , which consists of a set of  $k$  sub-objects. Suppose all clients in  $\beta$  behave in an well-formed manner. Suppose  $\beta$  contains no incomplete transactions and let  $\Pi$  be the set of transactions in  $\beta$ . Suppose there exists an irreflexive partial ordering ( $\prec$ ) among the transactions in  $\Pi$ , such that,*

- P1 *For any transaction  $\pi \in \Pi$  there are only a finite number of transactions  $\phi \in \Pi$  such that  $\phi \prec \pi$ ;*
- P2 *If the response event for  $\pi$  precedes the invocation event for  $\phi$  in  $\beta$ , then it cannot be that  $\phi \prec \pi$ ;*
- P3 *If  $\pi$  is a WRITE transaction in  $\Pi$  and  $\phi$  is any transaction in  $\Pi$ , then either  $\pi \prec \phi$  or  $\phi \prec \pi$ ; and*
- P4 *A tuple  $\mathbf{v} \equiv (v_{i_1}, v_{i_2}, \dots, v_{i_q})$  returned by a  $READ(o_{i_1}, o_{i_2}, \dots, o_{i_q})$ , where  $q$  is any positive integer,  $1 \leq q \leq k$ , is such that  $v_{i_j}$   $j \in \{1, \dots, q\}$  is written in  $\beta$  by the last preceding (w.r.t.  $\prec$ ) WRITE transaction that contains a  $write(o_{i_j}, *)$ , or the initial value  $v_{i_j}^0$  if no such WRITE exists in  $\beta$ .*

*Then execution  $\beta$  is strictly serializable.*

**Proof.** We discuss how to insert a serialization point  $*_\pi$  in  $\beta$  for every transaction  $\pi \in \Pi$ . First, we add  $*_\pi$  immediately after the latest of the invocations of  $\pi$  or  $\phi \in \Pi$  such that  $\phi \prec \pi$ . We want to stress that any complete operation  $\pi$  in  $\Pi$  refers to an operation, with the invocation and response events, whereas  $\pi_*$  refers to a point in the executions. Note that according to condition *P1* for  $\pi$  there are only finite number of such invocations in  $\beta$ , therefore,  $\pi_*$  is well-defined for any  $\pi \in \Pi$ . Now, since the order of the invocation events of the transactions in  $\Pi$  are already defined, therefore, the order of the corresponding set of serialization points are well-defined, except for the case when more than one serialization points are placed immediately after an invocation. In the case such multiple serialization points corresponding to an invocation we order these serialization points in accordance with the  $\prec$  relation of the underlying transactions.

Next, we show that for any pair of transactions  $\phi, \pi \in \Pi$  if  $\phi \prec \pi$  then in  $\beta$  we have  $*_\phi$  precedes  $*_\pi$ . Suppose  $\phi \prec \pi$ . By construction, each of  $\pi_*$  and  $\phi_*$  appear immediately after some invocation, in  $\beta$ , of some transaction in  $\Pi$ . If both  $\pi_*$  and  $\phi_*$  appear immediately after the same invocation, then since  $\phi \prec \pi$ , by construction of  $\pi_*$ , the serialization point  $\pi_*$  appears in  $\beta$  after  $\phi_*$ . Also, if the invocations after which  $\pi_*$  and  $\phi_*$  appear are distinct, then by construction of  $\pi_*$  the serialization point  $\pi_*$  appear after  $\phi_*$  in  $\beta$  since  $\phi \prec \pi$ .

Next we argue that each  $*_\pi$  serialization point for any  $\pi \in \Pi$  is placed between the invocation  $INV(\pi)$  and responses  $RESP(\pi)$ . By construction,  $*_\pi$  is after  $INV(\pi)$  in  $\beta$ . To show that  $*_\pi$  is before  $RESP(\pi)$  for the sake of contradiction assume that  $*_\pi$  appears after  $RESP(\pi)$ . By construction,  $*_\pi$  must be after  $INV(\psi)$  for some  $\psi \in \Pi$  and  $\psi \neq \pi$ , then by the condition of construction of  $\pi_*$  we have  $\psi \prec \pi$ . But from above  $INV(\psi)$  occurs after  $RESP(\pi)$ , i.e.,  $\pi$  completes before  $\psi$  is invoked which means, by property *P2*, we cannot have  $\psi \prec \pi$ , which is a contradiction.

Next, we show that if we were to shrink the transactions intervals to their corresponding serialization points, the resulting trace would be a trace of the underlying data type  $\mathcal{O}_T$ . In other words, we show any  $READ(o_{i_1}, o_{i_2}, \dots, o_{i_q})$  returns the values  $(v_{i_1}, v_{i_2}, \dots, v_{i_q})$ , such that each value  $v_{i_j}$ ,  $j \in [q]$ , was written by the immediately preceding (w.r.t. the serialization points)  $WRITE$  that contained  $write(o_{i_j}, v_{i_j})$  or the initial values if no such previous  $WRITE$  exists. Let us denote the set of  $WRITE$ s that precedes (w.r.t.  $\prec$ )  $\pi$  by  $\Pi_W^{\prec \pi}$ , i.e.,  $\phi \in \Pi_W^{\prec \pi}$   $\phi$  is a write and  $\phi \prec \pi$ . By property *P3*, all transactions in  $\Pi_W^{\prec \pi}$  are totally-ordered. By property *P4*,  $v_{i_j}$  must be the value updated by the most recent  $WRITE$  in  $\Pi_W^{\prec \pi}$ . Since the total order of serialization points are consistent with  $\prec$  and hence the  $v_{i_j}$  corresponds to the write operation of a  $WRITE$  transaction with the most recent serialization point and contains a operation of type  $write(o_{i_j}, *)$ . ◀

### E.3 SNOW with C2C Communication

In this section, we show that SNOW is possible in the *multiple-writers single-reader* (MWSR) setting when client-to-client communication is allowed. In particular, we present an algorithm  $A$ , which has all SNOW properties in such setting. We consider a system that has  $\ell \geq 1$  writers with ids  $w_1, w_2 \dots w_\ell \in \mathcal{W}$ , one reader  $r$ , and  $k \geq 1$  servers with ids  $s_1, s_2 \dots s_k \in \mathcal{S}$ . Client-to-client communication is allowed. The pseudocode for algorithm  $A$  is presented in Pseudocode 4. We use keys to uniquely identify a  $WRITE$  transaction in algorithm  $A$ . A key  $\kappa \in \mathcal{K}$  is defined as a pair  $(z, w)$ , where  $z \in \mathbb{N}$ , and  $w \in \mathcal{W}$  is the id of a writer.  $\mathcal{K}$  denotes the set of all possible keys. Also, with each transaction we associate a tag  $t \in \mathbb{N}$ .

**State variables:** (i) Any writer  $w$  has a counter  $z$  to keep track of the number of  $WRITE$  transactions it has invoked so far, initially 0. (ii) The reader  $r$  has an ordered list of elements,  $List$ , as  $(\kappa, (b_1, \dots, b_k))$ , where  $\kappa \in \mathcal{K}$  and  $(b_1, \dots, b_k) \in \{0, 1\}^k$ . Initially,



1074  $List = [(\kappa^0, (1, \dots, 1))]$ , where  $\kappa^0 \equiv (0, w_0)$ , and  $w_0$  is any place holder identifier for writer  
 1075 id.  $List$  can be thought of as an array, with 0 as the starting index. (iii) Each server  
 1076  $s_i \in \mathcal{S}$  has a set variable  $Vals$  with elements of key-value pairs  $(\kappa, v_i) \in \mathcal{K} \times \mathcal{V}_i$ . Initially,  
 1077  $Vals = \{(\kappa^0, v_i^0)\}$ .

1078 **Writer steps:** Any writer client,  $w \in \mathcal{W}$ , can issue a WRITE transaction that consists of a set of write operations by invoking the following operations  
 1079  $WRITE((o_{i_1}, v_{i_1}), (o_{i_2}, v_{i_2}), \dots, (o_{i_p}, v_{i_p}))$ , where  $I = \{i_1, i_2, \dots, i_p\}$  is any subset of  $p$  indices of  $[k]$ . We define the set  $S_I \triangleq \{s_{i_1}, s_{i_2}, \dots, s_{i_p}\}$ . This procedure consists of two  
 1080 consecutive phases: *write-value* and *inform-reader*. In the *write-value* phase,  $w$  creates a  
 1081 key  $\kappa$  as  $\kappa \equiv (z + 1, w)$ ; and also increments the local counter  $z$  by one. Then it sends  
 1082  $(WRITE-VALUE, (\kappa, v_i))$  to each server  $s_i$  in  $S_I$ , and awaits ACKNOWLEDGES from each server  
 1083 in  $S_I$ . After receiving all ACKNOWLEDGES,  $w$  initiates the *inform-reader* phase during which  
 1084 it sends  $(INFORM-READER, (\kappa, (b_1, \dots, b_k)))$  to  $r$ , where for any  $i \in [k]$ ,  $b_i$  is a boolean variable,  
 1085 such that  $b_i = 1$  if  $s_i \in S_I$ , otherwise  $b_i = 0$ . Essentially, such a  $(k + 1)$ -tuple identifies the  
 1086 set of objects that are updated during that WRITE transaction, i.e., if  $b_i = 1$  then object  
 1087  $o_i$  was updated during the execution of the WRITE transaction, otherwise  $b_i = 0$ . After  $w$   
 1088 receives ACKNOWLEDGE from  $r$  it completes the WRITE transaction.  
 1089

1090 **Reader steps:** We use the same notations for  $I$  and  $S_I$  as above for the set of  
 1091 indices and corresponding servers, possibly different across transactions. The procedure  
 1092  $READ(o_{i_1}, o_{i_2}, \dots, o_{i_p})$ , for any READ transaction, is initiated at reader  $r$ , where  
 1093  $o_{i_1}, o_{i_2}, \dots, o_{i_p}$  denotes the subset of objects  $r$  intends to read. This procedure consists of  
 1094 only one phase, *read-value*, of communication between the reader and the servers in  $S_I$ . Here  
 1095  $r$  sends the message  $(READ-VALUE, \kappa_i)$  to each server  $s_i \in S_I$ , where the  $\kappa_i$  is the key in the  
 1096 tuple  $(\kappa_i, (b_1, \dots, b_k))$  in  $List$  located at index  $j^*$  such that  $b_i = 1$  such that  $i \in I$ . After  
 1097 receiving the values  $v_{i_1}, v_{i_2}, \dots, v_{i_p}$  from all servers in  $S_I$ , where  $S_I \triangleq \{s_{i_1}, s_{i_2}, \dots, s_{i_p}\}$ , the  
 1098 transaction completes by returning  $(v_{i_1}, \dots, v_{i_p})$ .  
 1099

1100 If reader  $r$  receives a message  $(INFORM-READER, (\kappa, (b_1, \dots, b_k)))$  from any writer  $w$ , then  
 1101  $r$  appends  $(\kappa, (b_1, \dots, b_k))$  to its  $List$ , and responds to  $w$  with ACKNOWLEDGE and  $t_w = |List|$ ,  
 1102 i.e., number of elements in  $List$ . The order of the elements in  $List$  corresponds to the order  
 1103 the WRITE transactions, the order of the incoming INFORM-READER updates, as seen by  
 1104 the reader.

1105 **Server steps:** The server protocol consists of two procedures corresponding to the  
 1106 messages containing the tags WRITE-VALUE and READ-VALUE. The first procedure is used if  
 1107 a server  $s_i$  receives a message  $(WRITE-VALUE, (\kappa, v_i))$  from a writer  $w$ , it adds  $(\kappa, v_i)$  to its  
 1108 set variable  $Vals$  and sends ACKNOWLEDGE back to  $w$ . The second procedure is used if  $s_i$   
 1109 receives a message, i.e.,  $(READ-VALUE, \kappa_i)$ , from  $r$ , then it responds with  $v_i$  such that  $(\kappa_i, v_i)$   
 1110 is in its  $Vals$ .

1111 The following theorem proves that algorithm  $A$  respects all SNOW properties. Due to  
 1112 space constraints, the proof of the theorem can be found in Appendix E.3.

1113 **► Theorem 32.** *Any well-formed and fair execution of algorithm  $A$  is a wait-free implemen-*  
 1114 *tation of transaction processing in the MWSR setting with for objects of type  $\mathcal{O}_T$ , consisting*  
 1115 *of objects  $o_1, o_2, \dots, o_k$  maintained by the servers  $s_1, s_2, \dots, s_k$ , respectively; and it respects*  
 1116 *the SNOW properties and WRITES transactions are live.*

1117 **Proof.** Below we show that  $A$  satisfies the SNOW properties.

1118 **S property:** Let  $\beta$  be any fair execution of  $A$  and suppose all clients in  $\beta$  behave in a well-  
 1119 formed manner. Suppose  $\beta$  contains no incomplete transactions and let  $\Pi$  be the set of  
 1120 transactions in  $\beta$ . We define an irreflexive partial ordering  $(\prec)$  among the transactions in  $\Pi$



■ **Algorithm 4** The protocol for a writer  $w$ , reader  $r$  and server  $s_i$  for algorithm  $A$ .

<b>At writer <math>w</math></b>		Await ACKNOWLEDGE from $s_i$ for every $i \in I$ .	
<i>State Variables at <math>w</math>:</i>			
$z \in \mathbb{N}$ , initially 0		8: <u>inform-reader:</u>	
		<b>for</b> $i \in [k]$ <b>do</b>	
<b>WRITE</b> $((o_{i_1}, v_{i_1}), (o_{i_2}, v_{i_2}), \dots, (o_{i_p}, v_{i_p}))$		10: <b>if</b> $i \in I$ <b>then</b>	
2: <u>write-value:</u>		$b_i \leftarrow 1$	
$\kappa \leftarrow (z + 1, w); z \leftarrow z + 1$		12: <b>else</b>	
4: $I \triangleq \{i_1, i_2, \dots, i_p\}$		$b_i \leftarrow 0$	
<b>for</b> $i \in I$ <b>do</b>		14:     Send (INFORM-READER, $(\kappa, (b_1, \dots, b_k))$ ) to $r$	
6:     Send WRITE-VALUE, $(\kappa, v_{s_i})$ to $s_i$		Receive (ACKNOWLEDGE, $t_w$ ) from $r$	
16: _____			
<b>At reader <math>r</math></b>		Await responses $v_i$ from $s_i$ for each $i \in I$	
<i>State Variables at <math>r</math>:</i>		$t_r \triangleq \max_{1 \leq j \leq  List } \{j : List[j].b_i = 1 \wedge i \in I\}$	
$List$ , a list of elements in $\mathcal{K} \times \{0, 1\}^k$ ,		26:     Return $(v_{i_1}, v_{i_2}, \dots, v_{i_p})$	
18:     initially $[(\kappa^0, 1, \dots, 1)]$		_____	
<b>READ</b> $(o_{i_1}, o_{i_2}, \dots, o_{i_p})$		<b>Response routines</b>	
<u>read-value:</u>		28: <b>On recv</b> (INFORM-READER, $(\kappa, (b_1, \dots, b_k))$ ) <b>from</b> $w$ :	
20: $I \triangleq \{i_1, i_2, \dots, i_p\}$		$List \leftarrow List \oplus (\kappa, (b_1, \dots, b_k))$ // $\oplus$ for ap-	
<b>for</b> $i \in I$ <b>do</b>		pend	
22: $j^* \leftarrow \max_{1 \leq j \leq  List } \{j : List[j].b_i = 1\}$		30: $tag \leftarrow  List $ // $ \cdot $ size of the list	
$\kappa_i = List[j^*].\kappa$		Send (ACKNOWLEDGE, $tag$ ) to $w$	
24:     Send (READ-VALUE, $\kappa_i$ ) to $s_i$		32: _____	
<b>At server <math>s_i</math> for any <math>i \in [k]</math></b>		$Vals \leftarrow Vals \cup \{(\kappa, v)\}$	
<i>State Variables:</i>		Send ACKNOWLEDGE to writer $w$ .	
$Vals \subset \mathcal{K} \times \mathcal{V}_i$ , initially $\{(t_{key}^0, v_i^0)\}$		36: <b>On recv</b> (READ-VALUE, $\kappa$ ) <b>from</b> reader $r$ :	
34: <b>On recv</b> (WRITE-VALUE, $(\kappa, v)$ ) <b>from</b> writer $w$ :		Send $v$ s.t. $(\kappa, v) \in Vals$ to reader $r$	

1121 as follows: if  $\phi$  and  $\pi$  are any two distinct transactions in  $\Pi$  then we say  $\phi \prec \pi$  if either (i)  
 1122  $tag(\phi) < tag(\pi)$  or (ii)  $tag(\phi) = tag(\pi)$  and  $\phi$  is a WRITE and  $\pi$  is a READ. We will prove  
 1123 the  $S$  (strict-serializability) property of  $A$  by proving that the properties  $P1$ ,  $P2$ ,  $P3$  and  
 1124  $P4$  of Lemma ?? hold for  $\beta$ .

1125  **$P1$ :** If  $\pi$  is a READ then since all READS are invoked by a single reader  $r$  and in a  
 1126 well-formed manner, therefore, there cannot be an infinite number of READS such that they  
 1127 all precede  $\pi$  (w.r.t.  $\prec$ ). Now, suppose  $\pi$  is a WRITE. Clearly, from an inspection of the  
 1128 algorithm,  $tag(\pi) \in \mathbb{N}$ . From inspection of the algorithm, each WRITE increases the size of  
 1129  $List$ , and the value of the tags are defined by the size of  $List$ . Therefore, there can be at  
 1130 most a finite number of WRITES such that can precede  $\pi$  (w.r.t.  $\prec$ ) in  $\beta$ .

1131  **$P2$ :** Suppose  $\phi$  and  $\pi$  are any two transactions in  $\Pi$ , such that,  $\pi$  begins after  $\phi$  completes.  
 1132 Then we show that we cannot have  $\pi \prec \phi$ . Now, we consider four cases, depending on  
 1133 whether  $\phi$  and  $\pi$  are READS or WRITES.

1134 (a)  $\phi$  and  $\pi$  are WRITES invoked by writers  $w_\phi$  and  $w_\pi$ , respectively. Since the size of  $List$ ,  
 1135 in  $r$ , grows monotonically with each WRITE hence  $w_\pi$  receives the tag at least as high as  
 1136  $tag(\phi)$ , so  $\pi \not\prec \phi$ .

- 1137 (b)  $\phi$  is a WRITE,  $\pi$  is a READ transactions invoked by writer  $w_\phi$  and  $r$ , respectively. Since  
 1138 the size of *List*, in  $r$ , grows monotonically, and because  $w_\pi$  invokes  $\pi$  after  $\phi$  completes  
 1139 hence  $\text{tag}(\pi)$  is at least as high as  $\text{tag}(\phi)$ , so  $\pi \not\prec \phi$ .
- 1140 (c)  $\phi$  and  $\pi$  are READS invoked by reader  $r$ . Since the size of *List*, in  $r$ , grows monotonically,  
 1141 hence  $w_\pi$  invoked  $\pi$  after  $\phi$  completes hence  $\text{tag}(\pi)$  is at least as high as  $\text{tag}(\phi)$ , so  $\pi \not\prec \phi$ .
- 1142 (d)  $\phi$  is a READ,  $\pi$  is a WRITE invoked by reader  $r$  and  $w_\pi$ , respectively. This case is simple  
 1143 because new values are added to *List* only by writers, and  $\text{tag}(\pi)$  is larger than the tag  
 1144 of  $\phi$  and hence  $\pi \not\prec \phi$ .

1145 *P3*: This is clear by the fact that any WRITE transaction always creates a unique tag and  
 1146 all tags are totally ordered since they all belong to  $\mathbb{N}$

1147 *P4*: Consider a READ  $\rho$  as  $\text{READ}(o_{i_1}, o_{i_2}, \dots, o_{i_q})$ , in  $\beta$ . Let the returned value from  $\rho$   
 1148 be  $\mathbf{v} \equiv (v_{i_1}, v_{i_2}, \dots, v_{i_q})$  such that  $1 \leq i_1 < i_2 < \dots < i_q \leq k$ , where value  $v_{i_j}$  corresponds to  
 1149  $o_{i_j}$ . Suppose  $\text{tag}(\rho) \in \mathbb{N}$  was created during some WRITE transaction, say  $\phi$ , i.e.,  $\phi$  is the  
 1150 WRITE that added the elements in index  $(\text{tag}(\rho) - 1)$  of *List*. Note that element in index 0  
 1151 contains the initial value. Now we consider two cases:

1152 *Case  $\text{tag}(\rho) = 1$* . We know that it corresponds the initial default value  $v_i^0$  at each  
 1153 sub-object  $o_i$ , and this equates to  $\rho$  returning the default initial value for each sub-object.

1154 *Case  $\text{tag}(\rho) > 1$* . Then we argue that there exists no WRITE transaction, say  $\pi$ , that  
 1155 updated object  $o_{i_j}$ , in  $\beta$ , such that,  $\pi \neq \phi$  and  $\rho$  returns values written by  $\pi$  and  $\phi \prec \pi \prec \rho$ .  
 1156 Suppose we assume the contrary, which means  $\text{tag}(\phi) < \text{tag}(\pi) < \text{tag}(\rho)$ . The latter implies  
 1157  $\text{tag}(\phi) = \text{tag}(\pi)$  which is not possible because this contradicts the fact that for any two  
 1158 distinct WRITES  $\text{tag}(\phi) \neq \text{tag}(\pi)$  in any execution of  $A$ .

1159 *N property*: By inspection of algorithm  $A$  for the response steps of the servers to the reader.

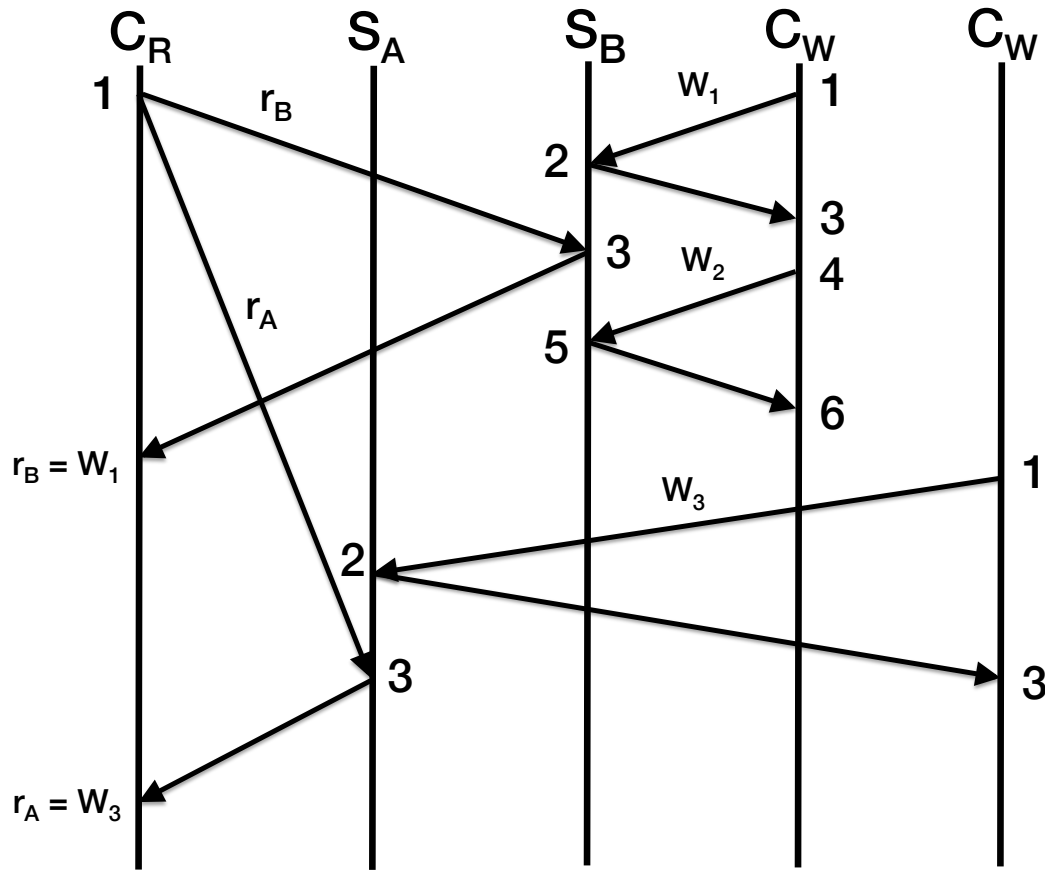
1160 *O property*: By inspection of the *read-value* phase: it consists of one round of communication  
 1161 between the reader and the servers, where the servers send only one version of the value of  
 1162 the object it maintains.

1163 *W property*: By inspection of the WRITE transaction steps, and and that writers always get  
 1164 to complete the transactions they invoke. ◀

## 1165 **F Eiger's READ transactions are not Strictly Serializable**

1166 Eiger [13] targeted a geo-replicated setting, which provided causal consistency across wide  
 1167 areas with fast READ and WRITE transactions handled by each datacenter locally. While  
 1168 not focusing on making transactions strictly serializable, Eiger made a claim that its READ  
 1169 transactions within the local datacenter were strictly serializable. Thus, in SNOW, there  
 1170 was a claim that Eiger was the only system, which had bounded latency while providing  
 1171 the strongest guarantees, and the practical lower bound for property *O* in existing SNW  
 1172 algorithms was three because Eiger's READ transactions required at most three rounds,  
 1173 i.e., the fewest in all examined existing work, were non-blocking, compatible with WRITE  
 1174 transactions, and were claimed to be strictly serializable. In this section, we correct this  
 1175 claim to be that there was no algorithm, which had bounded latency while providing the  
 1176 strongest guarantees by proving that Eiger's READ transactions are not strictly serializable.  
 1177 In Section ??, we provide the first set of novel SNW algorithms that have bounded latency  
 1178 and prove they are correct.

1179 The insight on why Eiger's READ transactions are not strictly serializable is that Eiger  
 1180 uses logical timestamps, i.e., Lamport clocks, to track the ordering of operations, and logical  
 1181 clocks are not able to identify the real-time ordering between operations that do not have



■ **Figure 3** An example execution that shows Eiger's READ transaction algorithm is not strictly serializable.  $w_1$ ,  $w_2$ , and  $w_3$  are three writes, where  $w_3$  is issued after  $w_2$  finishes.  $r_A$  and  $r_B$  are read operations from the same READ transaction  $R$ , which is concurrent with all three writes. Each number is a value of the logical clock on the machine (process) as a result of message exchange.

causal relationship, i.e., operations from different processes. Strict serializability, however, requires the real-time ordering to be respected.

Figure 3 is an example execution, which is allowed by Eiger but violates strict serializability. Real time goes downwards in the diagram. Each number is a value of the Lamport clock on the machine (process) as a result of message exchange. Initially all processes have Lamport clock value 0, and no messages happened before the execution in Figure 3. A READ transaction  $R = \{r_A, r_B\}$  reads values on  $S_A$  and  $S_B$  respectively. Due to the asynchronous nature of the network,  $r_B$  arrives on  $S_B$  before  $w_2$  and  $r_A$  arrives after  $w_3$ . Following the READ transaction algorithm of Eiger,  $r_A$  returns the value of  $w_3$  and its valid logical duration, i.e.,  $[2, 3]$ . Similarly,  $r_B$  returns the value of  $w_1$  and its valid logical duration, i.e.,  $[2, 3]$ . Because the two logical durations overlap, Eiger claims the combined values of  $r_A$  and  $r_B$  are consistent and accept them. However, because  $w_3$  starts after  $w_2$  finishes,  $w_3$  is in real time after  $w_2$ . By strict serializability, if a READ transaction sees the value of  $w_3$ , then it has to also observe the effect of  $w_2$ . Hence,  $R$ , which returns the values of  $w_1$  and  $w_3$ , violates strict serializability.

**G** SNW + One version, MWMR setting

■ **Algorithm 5** The protocol for any writer  $w$ , reader  $r$  or server  $s_i$  for algorithm  $B$ .

---

<b>At writer <math>w</math></b> <i>State Variables:</i> $z \in \mathbb{N}$ , initially 0 <b>WRITE</b> $((i_1, v_{i_1}), (i_2, v_{i_2}), \dots, (i_p, v_{i_p}))$ in $S_I$ . 2: $I \triangleq \{i_1, i_2, \dots, i_p\}$ <i>write-value:</i> 4: $\kappa \leftarrow (z + 1, w)$ ; $z \leftarrow z + 1$ <b>for</b> $i \in I$ <b>do</b>	6: Send WRITE-VALUE, $(\kappa, v_{s_i})$ to server $s_i$ 12: <b>await</b> ACKNOWLEDGE from server $s_i$ 14: Send (UPDATE-COORD, $(\kappa, (b_1, \dots, b_k)))$ to $s^*$ Receive (ACKNOWLEDGE, $t_w$ ) from coordinator $s^*$ 8: <i>update-coord:</i> <b>for</b> $i \in [k]$ <b>do</b> 10: <b>if</b> $i \in I$ <b>then</b>	$b_i \leftarrow 1$ <b>else</b> $b_i \leftarrow 0$ 14: Send (UPDATE-COORD, $(\kappa, (b_1, \dots, b_k)))$ to $s^*$ Receive (ACKNOWLEDGE, $t_w$ ) from coordinator $s^*$
--	---	--

---

16: <b>At reader <math>r</math></b> <b>READ</b> $(o_{i_1}, o_{i_2}, \dots, o_{i_p})$ $I \triangleq \{i_1, i_2, \dots, i_p\}$ 18: <i>get-tag-array:</i> Send (GET-TAG-ARRAY) to server $s^*$ 22: <i>read-value:</i> <b>for</b> $i \in I$ <b>do</b>	20: Receive response $(t_r, (\kappa_1, \kappa_2, \dots, \kappa_k))$ from $s^*$ 24: Wait responses as $v_i$ for each $s_i \in S$ Return $(v_{i_1}, v_{i_2}, \dots, v_{i_p})$	Send (READ-VALUE, $\kappa_i$ ) to $s_i$ 24: Wait responses as $v_i$ for each $s_i \in S$ Return $(v_{i_1}, v_{i_2}, \dots, v_{i_p})$
--	---	--

---

26: <b>At server <math>s_i</math> for any <math>i \in [k]</math></b> <i>State Variables:</i> $Vals \subset \mathcal{K} \times \mathcal{V}_i$ , initially $\{(\kappa^0, v_i^0)\}$ $List$ , a list of $\mathcal{K} \times \{0, 1\}^k$ , initially $[(\kappa^0, (1, \dots, 1))]$	<b>On recv</b> (UPDATE-COORD, $(\kappa, (b_1, \dots, b_k)))$ <b>from</b> $w$ : 32: $List \leftarrow List \oplus (\kappa, (b_1, \dots, b_k))$ /* used only by $s^*$ */ // $\oplus$ for append $tag \leftarrow  List  //  \cdot $ size of the list 34: Send (ACKNOWLEDGE, $tag$ ) to $w$ 28: <b>On recv</b> (WRITE-VALUE, $(\kappa, v)$ ) <b>from</b> $w$ : $Vals \leftarrow Vals \cup \{(\kappa, v)\}$ 30: Send ACKNOWLEDGE to writer $w$ . 36: <b>On recv</b> (READ-VALUE, $\kappa$ ) <b>from</b> $r$ : Send $v_i$ s.t. $(\kappa, v) \in Vals$ to $r$	<b>On recv</b> GET-TAG-ARRAY <b>from</b> $r$ : 38: <b>for</b> $i \in [k]$ <b>do</b> $j^* \leftarrow \max\{j : List[j].b_i =$ $1\}$ $\kappa_i = List[j^*].\kappa$ $t_r \triangleq \max_{1 \leq j \leq  List } \{j : List[j].b_i = 1 \wedge i \in I\}$ 42: Send $(t_r, (\kappa_1, \kappa_2, \dots, \kappa_k))$ to $r$
--	--	--

---

► **Theorem 33.** Any well-formed and fair execution of algorithm  $B$  is an implementation of an object of type  $\tilde{O}_T$  in the MWMR setting, with no client-to-client communication, comprising of objects  $o_1, o_2, \dots, o_k$  stored in servers  $s_1, s_2, \dots, s_k$ , respectively; and it satisfies the SNoW properties.

**Proof.** Below we show that algorithm  $B$  satisfies the SNoW properties.  
 $S$  property: Let  $\beta$  be any fair execution of  $B$  and suppose all clients in  $\beta$  behave in a well-formed manner. Suppose  $\beta$  contains no incomplete transactions and let  $\Pi$  be the set of transactions in  $\beta$ . We define an irreflexive partial ordering ( $\prec$ ) in  $\Pi$  as follows: if  $\phi$  and  $\pi$  are any two distinct transactions in  $\Pi$  then we say  $\phi \prec \pi$  if either (i)  $tag(\phi) < tag(\pi)$  or (ii)  $tag(\phi) = tag(\pi)$  and  $\phi$  is a WRITE transaction and  $\pi$  is a READ transaction. Below we prove the  $S$  property of  $B$  by showing that properties  $P1$ ,  $P2$ ,  $P3$  and  $P4$  of Lemma ?? hold for  $\beta$ .

$P1$ : Clearly, from an inspection of the algorithm,  $tag(\pi) \in \mathbb{N}$ . From inspection of the algorithm, each WRITE transaction increases the size of  $List$ , and the value of the tags are defined by the size of  $List$ . Therefore, there can be at most a finite number of WRITE transactions such that can precede  $\pi$  (w.r.t.  $\prec$ ) in  $\beta$ . On the other hand, if  $\pi$  is a READ transaction then since all READ transactions are invoked by readers in a well-formed manner, and there are only finite number of readers therefore, there cannot be an infinite number of READ transactions such that they all precede  $\pi$  (w.r.t.  $\prec$ ).

$P2$ : Suppose  $\phi$  and  $\pi$  are any two transactions in  $\Pi$ , such that,  $\pi$  begins after  $\phi$  completes.

1218 Then we show that we have we cannot have  $\pi \prec \phi$ . Now, we consider four cases, depending  
 1219 on whether  $\phi$  and  $\pi$  are READ transactions or WRITE transactions.

- 1220 (a)  $\pi$  and  $\phi$  are WRITE transactions invoked by writers  $w_\pi$  and  $w_\phi$ , respectively. Since the  
 1221 size of  $List$ , in  $s^*$  grows monotonically due to each WRITE transaction hence  $w_\pi$  receives  
 1222 the tag from  $s^*$  at least as high as  $tag(\phi)$ , so  $\pi \not\prec \phi$ .  
 1223 (b)  $\pi$  is a READ transaction,  $\phi$  is a WRITE transaction invoked by reader  $r_\pi$  and writer  
 1224  $w_\phi$ , respectively. Since the size of  $List$ , in  $s^*$ , grows monotonically, because  $r_\pi$  invokes  $\pi$   
 1225 after  $\phi$  completes hence  $tag(\pi)$  is at least as high as  $tag(\phi)$ , so  $\pi \not\prec \phi$ .  
 1226 (c)  $\pi$  and  $\phi$  are both READ transactions invoked by readers  $r_\pi$  and  $r_\phi$ , respectively. Since  
 1227 the size of  $List$ , in  $s^*$ , grows monotonically, because  $w_\pi$  invokes  $\pi$  after  $\phi$  completes hence  
 1228  $tag(\pi)$  is at least as high as  $tag(\phi)$ , so  $\pi \not\prec \phi$ .  
 1229 (d)  $\pi$  is a WRITE transaction,  $\phi$  is a READ transaction invoked by writer  $w_\pi$  and reader  
 1230  $r_\phi$ , respectively. This case is simple because new values are added to  $List$ , in  $s^*$ , only by  
 1231 writers, and  $tag(\pi)$  has to be larger than the tag of  $\phi$  and hence  $\pi \not\prec \phi$ .

1232 *P3*: This is from the fact that any WRITE transaction always creates a unique tag and  
 1233 all tags are totally ordered since they all belong to  $\mathbb{N}$

1234 *P4*: Consider a READ transaction  $\rho$  as  $READ(o_{i_1}, o_{i_2}, \dots, o_{i_q})$ , in  $\beta$ . Let the returned  
 1235 value from  $\rho$  be  $\mathbf{v} \equiv (v_{i_1}, v_{i_2}, \dots, v_{i_q})$  such that  $1 \leq i_1 < i_2 < \dots < i_q \leq k$ , where value  $v_{i_j}$   
 1236 corresponds to  $o_{i_j}$ . Suppose  $tag(\rho) \in \mathbb{N}$  was created during some WRITE transaction, say  $\phi$ ,  
 1237 i.e.,  $\phi$  is the WRITE transaction that added the elements in index  $(tag(\rho) - 1)$  of  $List$  at  
 1238 the coordinator  $s^*$ . Note that element in index 0 contains the initial value. Now we consider  
 1239 two cases:

1240 *Case*  $tag(\rho) = 1$ . We know that it corresponds the initial default value  $v_i^0$  at each  
 1241 sub-object  $o_i$ , and this equates to  $\rho$  returning the default initial value for each sub-object.

1242 *Case*  $tag(\rho) > 1$ . Then we argue that there exists no WRITE transaction, say  $\pi$ , that  
 1243 updated object  $o_{i_j}$ , in  $\beta$ , such that,  $\pi \neq \phi$  and  $\rho$  returns values written by  $\pi$  and  $\phi \prec \pi \prec \rho$ .  
 1244 Suppose we assume the contrary, which means  $tag(\phi) < tag(\pi) < tag(\rho)$ . The latter implies  
 1245  $tag(\phi) = tag(\pi)$  which is not possible because this contradicts the fact that for any two  
 1246 distinct WRITE transactions  $tag(\phi) \neq tag(\pi)$  in any execution of  $B$ .

1247 *N, o and W properties*: Evident from an inspection of the algorithm. ◀

## 1248 **H SNW + One round, MWMR setting**

1249 ▶ **Theorem 34.** *Any well-formed and fair execution of algorithm  $B$  is an implementation*  
 1250 *of an object of type  $\tilde{O}_T$  in the MWMR setting, with no client-to-client communication,*  
 1251 *comprising of objects  $o_1, o_2, \dots, o_k$  stored in servers  $s_1, s_2, \dots, s_k$ , respectively; and it satisfies*  
 1252 *the  $SN\bar{o}W$  properties.*

1253 ▶ **Theorem 35.** *Any well-formed and fair execution of algorithm  $C$  is an implementation*  
 1254 *of an object of type  $\tilde{O}_T$  in the MWMR setting, with no client-to-client communication,*  
 1255 *comprising of objects  $o_1, o_2, \dots, o_k$  stored in servers  $s_1, s_2, \dots, s_k$ , respectively; and it satisfies*  
 1256 *the  $SN\bar{o}W$  properties.*

1257 **Proof.** Below we show that algorithm  $C$  satisfies the  $SN\bar{o}W$  properties.

1258 *S property*: Let  $\beta$  be any fair execution of  $B$  and suppose all clients in  $\beta$  behave in a well-  
 1259 formed manner. Suppose  $\beta$  contains no incomplete transactions and let  $\Pi$  be the set of  
 1260 transactions in  $\beta$ . We define an irreflexive partial ordering ( $\prec$ ) in  $\Pi$  as follows: if  $\phi$  and  $\pi$   
 1261 are any two distinct transactions in  $\Pi$  then we say  $\phi \prec \pi$  if either (i)  $tag(\phi) < tag(\pi)$  or (ii)

■ **Algorithm 6** The protocol for any writer  $w$ , reader  $r$  or server  $s_i$  for algorithm  $C$ .

---

<b>At writer <math>w</math></b>	6: Send WRITE-VALUE, $(\kappa, v_{s_i})$	$b_i \leftarrow 1$
<b>State Variables:</b>	to server $s_i$	<b>12: else</b>
$z \in \mathbb{N}$ , initially 0	Await ACKNOWLEDGE from	$b_i \leftarrow 0$
<b>WRITE</b> $((i_1, v_{i_1}), (i_2, v_{i_2}), \dots, (i_p, v_{i_p}))$ in $S_I$ .	<b>14: Send</b> (UPDATE-COORD,	
2: $I \triangleq \{i_1, i_2, \dots, i_p\}$	$(\kappa, (b_1, \dots, b_k))$ to $s^*$	
<u>write-value:</u>	Receive (ACKNOWLEDGE, $t_w$ )	
4: $\kappa \leftarrow (z + 1, w)$ ; $z \leftarrow z + 1$	from coordinator $s^*$	
<b>for <math>i \in I</math> do</b>	<b>8: update-coord:</b>	
	<b>for <math>i \in [k]</math> do</b>	
	<b>10: if <math>i \in I</math> then</b>	
16: <b>At reader <math>r</math></b>	Send (GET-TAG-ARRAY) to	$(t_r, (\kappa_1, \kappa_2, \dots, \kappa_k))$ from $s^*$
<b>READ</b> $(o_{i_1}, o_{i_2}, \dots, o_{i_p})$	server $s^*$	Wait responses as $Vals_i$ from
$I \triangleq \{i_1, i_2, \dots, i_p\}$	<b>20: for <math>i \in I</math> do</b>	each $s_i \in S_I$
18: <u>read-values-and-tags:</u>	Send (READ-VALUES) to $s_i$	<b>24: Return</b> $(v_{i_1}, v_{i_2}, \dots, v_{i_p})$ s.t.
	<b>22: Receive</b> response	$(\kappa_j, v_j) \in Vals_j, j \in I$
<hr/>		
<b>At server <math>s_i</math> for any <math>i \in [k]</math></b>		
26: <b>State Variables:</b>	<b>30: On recv</b> (UPDATE-COORD, $(\kappa, (b_1, \dots, b_k^*))$ ) <b>from <math>s^*</math>:</b>	<b>36: On recv</b> GET-TAG-ARRAY <b>from <math>r</math>:</b>
$Vals \subset \mathcal{K} \times \mathcal{V}_i$ , initially	$List \leftarrow List \oplus (\kappa, (b_1, \dots, b_k^*))$	<b>for <math>i \in [k]</math> do</b>
$\{(\kappa^0, v_i^0)\}$	// $\oplus$ for append	<b>38: <math>j^* \leftarrow \max\{j : List[j].b_i =</math></b>
$List$ , a list of $\mathcal{K} \times \{0, 1\}^k$ , initially $[(\kappa^0, (1, \dots, 1))]$	<b>32: <math>tag \leftarrow  List  //  \cdot </math> size of the</b>	<b>1}</b>
	list	$\kappa_i = List[j^*].\kappa$
	Send (ACKNOWLEDGE, $tag$ ) to	<b>40: <math>t_r \triangleq \max_{1 \leq j \leq  List } \{j :</math></b>
<b>On recv</b> (WRITE-VALUE, $(\kappa, v)$ ) <b>from <math>w</math>:</b>		$List[j].b_i = 1 \wedge i \in I\}$
28: $Vals \leftarrow Vals \cup \{(\kappa, v)\}$	<b>34: On recv</b> (READ-VALUES) <b>from <math>r</math>:</b>	Send $(t_r, (\kappa_1, \kappa_2, \dots, \kappa_k))$ to
Send ACKNOWLEDGE to writer $w$ .	Send $Vals$ to $r$	$r$

---

1262  $tag(\phi) = tag(\pi)$  and  $\phi$  is a WRITE and  $\pi$  is a READ. Below we prove the  $S$  property of  $B$  by  
 1263 showing that properties  $P1, P2, P3$  and  $P4$  of Lemma ?? hold for  $\beta$ . The properties  $P1-P4$   
 1264 can be proved to hold in a manner very similar to algorithm  $B$  (Section ??). Therefore, we  
 1265 avoid repeating them.  
 1266  $N, \bar{o}$  and  $W$  properties: Evident from an inspection of the algorithm. ◀