



HAL
open science

A Formal Analysis of the Global Sequence Protocol

Hernán Melgratti, Christian Roldán

► **To cite this version:**

Hernán Melgratti, Christian Roldán. A Formal Analysis of the Global Sequence Protocol. 18th International Conference on Coordination Languages and Models (COORDINATION), Jun 2016, Heraklion, Greece. pp.175-191, 10.1007/978-3-319-39519-7_11 . hal-01631716

HAL Id: hal-01631716

<https://inria.hal.science/hal-01631716>

Submitted on 9 Nov 2017

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

A Formal Analysis of the Global Sequence Protocol ^{*}

Hernán Melgratti^{1,2} and Christian Roldán¹

¹ Departamento de Computación, FCEyN, Universidad de Buenos Aires.
² CONICET.

Abstract. The Global Sequence Protocol (GSP) is an operational model for replicated data stores, in which updates propagate asynchronously. We introduce the GSP-calculus as a formal model for GSP. We give a formal account for its proposed implementation, which addresses communication failures and compact representation of data, and use simulation to prove that the implementation is correct. Then, we use the GSP-calculus to reason about execution histories and prove ordering guarantees, such as read my writes, monotonic reads, causality and consistent prefix. We also prove that GSP extended with synchronous updates provides strong consistency guarantees.

1 Introduction

Cloud infrastructures provide data storages that are virtually unlimited, elastic (i.e., scalable at run-time), highly available and partition tolerant. This is achieved by replicating data over multiple servers. A client may perform update and read operations over any of these replicas and the store is responsible for keeping them synchronised. However, it is known (CAP theorem [7]) that any system cannot simultaneously provide availability, partition tolerance, and consistency. Thus, one of these properties has to be discarded. Today's popular data storages, such as Dynamo [6] and Cassandra [9], ensure availability and offer weaker notions of consistency, called *eventual consistency*. Roughly, eventual consistency guarantees that all updates will be delivered to the different replicas, which will eventually converge to the same state [1]. The storages adopt different strategies to achieve eventual consistency, which impact on the guarantees provided by the system, i.e., on the kind of inconsistencies or anomalies that are allowed to happen. For instance, a storage may resolve automatically conflicts introduced by concurrent updates (e.g., by using timestamps or causality) or may leave the problem to applications that read the database (like in Cassandra). In this way, the consistency model supported by a data store becomes crucial when writing applications.

Consequently, there has been a growing interest on establishing programming abstractions to help developers to deal with eventual consistent stores. For

^{*} Research partially supported by UBACyT project 2014-2017 20020130200092BA.

instance, commutative replicated data types [10] and cloud types [3] provide programmers with suitable data type abstractions that encapsulate issues related to eventual consistency. Recent proposals advocate declarative approaches for programming with eventual consistency, e.g., to automatically select the consistency level required from a store provided with a consistency contract for the application [11] or to prove that a given consistency level is adequate for preserving some data invariant [8]. With similar aims, the Global Sequence Protocol (GSP) [5] proposes an operational model to reason about applications running on top of replicated stores. Basically, the state of a store is represented as the sequence of updates that have led to it. Clients have their own copy of the state which they operate upon: each read and write operation has immediate effect over the local state and the system propagates changes to make all replicas consistent using a reliable total order broadcast protocol (RTOB). The RTOB protocol guarantees that all messages are delivered in the same total order to all clients. Replicas rely on the order generated by RTOB to converge to the same state. In the very basic model, called core GSP, each client interacts with its local state by performing read and write operations. Albeit simple, this model introduces some subtleties when programming because it does not ensure read stability (i.e., two successive reads may return different values) nor atomicity of several updates (i.e., another client may partially observe the effects of a sequence of updates). To overcome these limitations, three synchronisation primitives, namely `pull`, `push` and `confirmed`, allow programmers to control the propagation of changes. It has been shown that this model can be implemented so to handle communication failures and to represent updates efficiently by using two type of objects: states and deltas. Both models, i.e. the idealised one and its implementation, have been defined in terms of a reference implementation.

In this paper, we propose a formal account for each model: the GSP and IGSP calculi (§ 2 and § 3). We prove that the behaviour of a program running over IGSP can be observed over the idealised model. Technically, we show that each IGSP system can be simulated by the corresponding GSP system (§ 4). Then, we study and prove the consistency guarantees ensured by GSP. We rely on the characterisation of consistency guarantees in terms of abstract histories proposed in [2]. Abstract histories capture the visibility relation between actions and the arbitration order of updates in the system. Then, a wide-spectrum of consistency models can be characterised in terms of these two relations. In § 5, we show how to operationally associate abstract histories to concrete computations and prove that GSP enjoys properties such as *Monotonic Read*, *Causal Visibility* and *Consistent prefix*, among others. Finally, in § 6 we study the extension of GSP with synchronous write operations, which ensures strong consistency.

2 Global Sequence Protocol Calculus

2.1 Syntax

Clients interact with a store by performing operations in $\mathcal{U} \cup \mathcal{R}$: an element in \mathcal{U} denotes an update operation, while one in \mathcal{R} stands for a read operation.

(NATURALS)	$j, k, n \in \mathbb{N}$		
(UPDATE)	$\mathcal{U} = \{u, u', \dots, u_0, \dots\}$	(UPD SEQ)	$\mathbf{u} ::= \epsilon \mid u^{\mathfrak{v}} \cdot \mathbf{u}$
(READ)	$\mathcal{R} = \{r, r', \dots, r_0, \dots\}$	(BLOCK SEQ)	$\mathbf{b} ::= \epsilon \mid (\mathbf{u}_{\top}) \cdot \mathbf{b}$
(EVENT)	$\mathbb{V} = \{\mathfrak{v}, \mathfrak{v}', \dots, \mathfrak{v}_0, \dots\}$	(SYSTEM)	$N ::= S \parallel C$
(VAR)	$\mathcal{X} = \{x, x', \dots, x_0, \dots\}$	(STORE)	$S ::= \mathbf{b}$
(IDS)	$\mathcal{I} = \{i, j, i', \dots, i_0, \dots\}$	(CLIENT)	$C ::= 0 \mid \langle P, \mathbf{u}_{\top}, \mathbf{b}, \mathbf{b}, k, j \rangle_i \mid C \parallel C$
(PROGRAM)	$P ::= \text{update}(u); P \mid \text{let } x = \text{read}(r) \text{ in } P \mid \text{pull}; P \mid \text{push}; P \mid \text{let } x = \text{confirmed in } P$		

Fig. 1. Syntax of the GSP calculus

No operation can simultaneously read and update a store, therefore we assume $\mathcal{U} \cap \mathcal{R} = \emptyset$. We write u, u', u'', \dots for updates and r, r', r'', \dots for reads.

The state of a store is represented by a sequence of updates. For technical convenience (particularly in § 5), we distinguish different executions of the same operation. Formally, stores associate each update with a fresh event identifier. We assume a set \mathbb{V} of event identifiers $\mathfrak{v}, \mathfrak{v}_0, \dots, \mathfrak{v}', \dots$ and write $u^{\mathfrak{v}}$ for the update u associated with the event \mathfrak{v} .

We use \mathbf{u} to denote sequences of decorated updates and (\mathbf{u}) for an atomic block of updates. We write \mathbf{b} for a sequence of blocks. We denote the empty sequence with ϵ and use the usual operations on sequences such as $\mathbf{b}[i]$ to denote the i -th element of \mathbf{b} , $\mathbf{b}[i..j]$ for the subsequence of \mathbf{b} from position i to j , $|\mathbf{b}|$ for its length and $\mathbf{b} \setminus \mathbf{b}'$ for the relative complement of \mathbf{b} in \mathbf{b}' . Additionally, $\underline{\mathbf{b}}$ stands for the plain sequence of updates in \mathbf{b} (i.e., without any separation in blocks).

We rely on the countable sets \mathcal{X} of program variables x, x', \dots and \mathcal{I} of client identifiers i, i', \dots, i_1, \dots .

Definition 2.1 (GSP Language). *The set of GSP terms is given by the grammar in Fig. 1.*

A GSP system N consists of a store and zero or more clients. The global store S is completely defined by its state, which consists of a sequence of blocks. The term $\langle P, \mathbf{u}_{\top}, \mathbf{b}_{\mathcal{S}}, \mathbf{b}_{\mathcal{P}}, k, j \rangle_i$ stands for a client identified by i and engaged on the execution of the program P . The remaining elements are used to describe the state of the local replica: \mathbf{u}_{\top} contains the updates that have been made locally and are part of an unfinished block; $\mathbf{b}_{\mathcal{S}}$ models the communication buffer, which keeps all blocks completed by the client but not received by the global store; $\mathbf{b}_{\mathcal{P}}$ is the pending buffer, which contains all completed blocks that are unconfirmed by the global store. For simplicity, we do not have an explicit replica of the global store in each client; we use instead a natural number k to indicate the portion of the global state that is known to the client. Specifically, the client i knows the sequence $S[0..k-1]$. Similarly, j indicates the number of updates received by the client that have not been added to the local replica, i.e., the client has received the updates contained in the segment $S[k..k+j-1]$.

A program P is built as a sequence of operations that interacts with the store: `read`(r), `update`(u), `pull`, `push`, `confirmed` (we postpone their description

until § 2.2). A program $\text{let } x = \dots \text{ in } P$ introduces a bound variable whose scope is P . The definition of free variables of a program is standard. We say that a process P is *closed* when it does not contain free variables. We keep the language for programs simple. We remark that this choice does not affect the results presented in this paper. Actually, we could just have characterised the behaviour of programs as a labelled transition system, but we prefer to have a syntax throughout the presentation.

Definition 2.2 (Well-formedness). A GSP system $N = C_0 \parallel \dots \parallel C_m \parallel S$ where $C_l = \langle P_l, \mathbf{u}_{\top l}, \mathbf{b}_{S_l}, \mathbf{b}_{P_l}, k_l, j_l \rangle_i$ for all $l \in \{0, \dots, m\}$ is well-formed if the following conditions hold

1. $i_l \neq i_{l'}$ for all $l \neq l'$;
2. $k_l + j_l \leq |S|$ for all l ;
3. $\mathbf{b}_{P_l} = \langle \mathbf{u}_1 \rangle \dots \langle \mathbf{u}_p \rangle \cdot \mathbf{b}_{S_l}$ and for all $1 \leq x < y \leq p$ there exists x', y' s.t. $S[x'] = \langle \mathbf{u}_x \rangle$, $S[y'] = \langle \mathbf{u}_y \rangle$ and $k_l \leq x' < y'$; and
4. $\mathbf{u} = S \cdot \mathbf{b}_{S_0} \dots \mathbf{b}_{S_m} \cdot \mathbf{u}_{\top 0} \dots \mathbf{u}_{\top m}$, if $\mathbf{u}[x] = \mathbf{u}'$, $\mathbf{u}[y] = \mathbf{u}'$ and $x \neq y$ then $\mathbf{v} \neq \mathbf{v}'$.

We require identifiers to univocally identify clients (1) and every local state to be consistent with the global store, i.e., a client can see at most every message in the store (2), all unconfirmed blocks in \mathbf{b}_{P_l} are either in the communication buffer \mathbf{b}_{S_l} or in the unseen part of the global store $\langle \mathbf{u}_1 \rangle \dots \langle \mathbf{u}_p \rangle$ (3). Moreover, an event identifier is associated with a unique update in the system (4). Hereafter, we assume every GSP system to be well-formed.

2.2 Operational Semantics

The operational semantics of GSP is given by a labelled transition system over well-formed terms, quotiented by the structural equivalence \equiv defined as the least equivalence such that \parallel is associative, commutative and has 0 as neutral element. The set of actions is given by the following grammar:

$$\lambda ::= \tau \mid rd(r) \mid wr(\mathbf{u}') \mid pull \mid push \mid cfm$$

As usual, τ stands for an internal, unobservable action, while the remaining ones correspond to the interaction of a client with the store. A label (λ, i) indicates that the client i performs the action λ . We write $\xrightarrow{\lambda}_i$ instead of $\xrightarrow{(\lambda, i)}$.

We now comment on the inference rules in Fig. 2. When a client performs an update (rule UPDATE), the change has only local effects: the sequence of local updates \mathbf{u}_{\top} is extended with the operation u decorated with a globally fresh identifier \mathbf{v} . We remark that decorations are used for technical reasons but they are operationally irrelevant (see § 5).

A client propagates its local changes to the global store by executing **push** (rule PUSH): all local changes in \mathbf{u}_{\top} will be transmitted as a block $\langle \mathbf{u}_{\top} \rangle$, i.e., as an atomic unit. Nevertheless, these changes are not made available immediately

$$\begin{array}{c}
\text{(UPDATE)} \\
\hline
\text{v fresh} \\
\hline
\langle \text{update}(u); P, \mathbf{u}_T, \mathbf{b}_S, \mathbf{b}_P, k, j \rangle_i \parallel N \xrightarrow{wr(u^\vee)}_i \langle P, \mathbf{u}_T \cdot u^\vee, \mathbf{b}_S, \mathbf{b}_P, k, j \rangle_i \parallel N \\
\text{(PUSH)} \\
\langle \text{push}; P, \mathbf{u}_T, \mathbf{b}_S, \mathbf{b}_P, k, j \rangle_i \parallel N \xrightarrow{push}_i \langle P, \epsilon, \mathbf{b}_S \cdot \langle \mathbf{u}_T \rangle, \mathbf{b}_P \cdot \langle \mathbf{u}_T \rangle, k, j \rangle_i \parallel N \\
\text{(SEND)} \\
\langle P, \mathbf{u}_T, \langle \mathbf{u}'_T \rangle \cdot \mathbf{b}_S, \mathbf{b}_P, k, j \rangle_i \parallel C \parallel S \xrightarrow{\tau}_i \langle P, \mathbf{u}_T, \mathbf{b}_S, \mathbf{b}_P, k, j \rangle_i \parallel C \parallel S \cdot \langle \mathbf{u}'_T \rangle \\
\text{(RECEIVE)} \\
\hline
k + j < |S| \\
\hline
\langle P, \mathbf{u}_T, \mathbf{b}_S, \mathbf{b}_P, k, j \rangle_i \parallel C \parallel S \xrightarrow{\tau}_i \langle P, \mathbf{u}_T, \mathbf{b}_S, \mathbf{b}_P, k, j + 1 \rangle_i \parallel C \parallel S \\
\text{(PULL)} \\
\langle \text{pull}; P, \mathbf{u}_T, \mathbf{b}_S, \mathbf{b}_P, k, j \rangle_i \parallel C \parallel S \xrightarrow{pull}_i \langle P, \mathbf{u}_T, \mathbf{b}_S, \mathbf{b}_P \setminus S[k..k + j - 1], k + j, 0 \rangle_i \parallel C \parallel S \\
\text{(READ)} \\
\hline
rvalue(r, S[0..k - 1] \cdot \mathbf{b}_P \cdot \mathbf{u}_T) = v \\
\hline
\langle \text{let } x = \text{read}(r) \text{ in } P, \mathbf{u}_T, \mathbf{b}_S, \mathbf{b}_P, k, j \rangle_i \parallel C \parallel S \xrightarrow{rd(r)}_i \langle P\{v/x\}, \mathbf{u}_T, \mathbf{b}_S, \mathbf{b}_P, k, j \rangle_i \parallel C \parallel S \\
\text{(CONFIRM)} \\
\hline
v = (\mathbf{b}_P \cdot \mathbf{u}_T == \epsilon) \\
\hline
\langle \text{let } x = \text{confirmed in } P, \mathbf{u}_T, \mathbf{b}_S, \mathbf{b}_P, k, j \rangle_i \parallel N \xrightarrow{cfm}_i \langle P\{v/x\}, \mathbf{u}_T, \mathbf{b}_S, \mathbf{b}_P, k, j \rangle_i \parallel N
\end{array}$$

Fig. 2. Operational semantics for GSP

at the global store because of the asynchronous communication model. In fact, the new block $\langle \mathbf{u}_T \rangle$ is added to the communication buffer \mathbf{b}_S , which contains all blocks that have not reached the global store. Also, $\langle \mathbf{u}_T \rangle$ is added to the pending messages \mathbf{b}_P . Rule SEND stands for a block that finally reaches the global store. Conversely, rule RECEIVE models the reception of a new update. The received update is not immediately added to the local replica. Actually, each client explicitly refreshes its local view by executing `pull` (rule PULL). At this time, all previously received updates j are incorporated to the local copy (i.e., k is changed to $k + j$). Additionally, all pending updates in the new fragment $S[k..k + j - 1]$ are removed from \mathbf{b}_P .

The semantics of operations is defined abstractly by the interpretation function $rvalue : \mathcal{R} \times \mathcal{U}^* \rightarrow \mathcal{V}$, i.e., a function that takes a read operation and a sequence of updates and returns a value in some domain \mathcal{V} . A read operation r is evaluated over the local state of the client (rule READ), i.e., the known prefix of the global store $S[0..k - 1]$ and the local updates in \mathbf{b}_P and \mathbf{u}_T . The value v is bound to the variable x , and hence all free occurrences of x in the continuation P are substituted by v . A client may perform `confirmed` to check whether its executed updates are already in the global store: this operation returns true only when the local buffers \mathbf{b}_P and \mathbf{u}_T are both empty (rule CONFIRM).

We remark that the operational semantics preserves well-formedness.

Lemma 2.1. *If N is well-formed and $N \xrightarrow{\lambda}_i N'$, then N' is well-formed.*

(STATE)	$\mathbf{s}, \mathbf{s}', \dots, \mathbf{s}_1, \dots \in State$	(SGMT SEQ)	$seg ::= \epsilon \mid seg \cdot seg$
(DELTA)	$\delta, \delta', \dots, \delta_1, \dots \in Delta$	(IN SRV)	$\mathbf{in}_s \in \mathcal{I} \rightarrow \mathbf{r}$
(ROUND)	$r ::= \langle i, n, \delta \rangle$	(OUT SRV)	$\mathbf{out}_s \in \mathcal{I} \rightarrow seg$
(RND SEQ)	$\mathbf{r} ::= \epsilon \mid r \cdot \mathbf{r}$	(SYSTEM)	$\mathbf{N} ::= \mathbf{S} \parallel \mathbf{C}$
(MAX RND)	$\mathbf{f}, \mathbf{f}', \dots, \mathbf{f}_1, \dots \in \mathcal{I} \rightarrow \mathbb{N}$	(SERVER)	$\mathbf{S} ::= \langle \mathbf{s}, \mathbf{f}, \mathbf{in}_s, \mathbf{out}_s \rangle$
(SEGMENT)	$seg ::= \langle \delta, \mathbf{f} \rangle \mid \langle \mathbf{s}, \mathbf{f} \rangle$	(CLIENTS)	$\mathbf{C} ::= 0 \mid \langle P, \mathbf{s}, \delta, \delta, n, \mathbf{r}, seg \rangle_i \mid \mathbf{C} \parallel \mathbf{C}$

Fig. 3. Syntax of the IGSP calculus

3 Implementation of GSP

The GSP model describes an idealised system that abstracts away from several implementation details, such as non-optimised representation of the state and unreliable communication. This section presents a formal model for the implementation proposed in [5].

3.1 Syntax

The implementation of GSP relies on a compact representation for states and updates. Their precise definition highly depends on the datatype of the values handled by the store, but they are characterised in terms of two abstract types: *State* and *Delta*, which provides the following operations [5]:

δ_\emptyset	: <i>Delta</i>	\emptyset	: <i>State</i>
<i>append</i>	: $Delta \times \mathcal{U} \rightarrow Delta$	<i>apply</i>	: $State \times Delta^* \rightarrow State$
<i>reduce</i>	: $Delta^* \rightarrow Delta$	<i>read</i>	: $\mathcal{R} \times State \rightarrow \mathcal{V}$

Constants δ_\emptyset and \emptyset denote the empty elements in their respective types. An object $\delta \in Delta$ describes the effects of a sequence of updates and is built by either appending an update to an existing delta (*append*) or combining together several deltas (*reduce*). Operation *read* is the interpretation function for operations (i.e., the implementation counterpart of function *rvalue*(-, -) used by the idealised model) and *apply* corresponds to state transformations.

Clients and the global store exchange δ objects to communicate changes. As each single δ may correspond to several update operations, clients send each δ accompanied by its own identifier and a sequence number n . Precisely, clients send rounds, i.e. triples $r = \langle i, n, \delta \rangle$. Differently, the global store sends segments $seg = \langle \delta, \mathbf{f} \rangle$, in which δ is accompanied by a function $\mathbf{f} \in \mathcal{I} \rightarrow \mathbb{N}$. In this way, the global store confirms all changes from client i until round $\mathbf{f}(i)$. To deal with crashes and recovery, the server may send segments of the form $\langle \mathbf{s}, \mathbf{f} \rangle$, which communicates a complete state instead of a delta object.

Definition 3.1 (GSP Language). *The set of IGSP terms is given by the grammar in Fig. 3.*

As for GSP, a system is composed by a global store \mathbf{S} and possibly many clients \mathbf{C} . A global store is modelled by a tuple $\langle \mathbf{s}, \mathbf{f}, \mathbf{in}_s, \mathbf{out}_s \rangle$ containing a

state \mathbf{s} , a function \mathbf{f} to keep track of processed rounds and the communication buffers \mathbf{in}_s and \mathbf{out}_s . There are two dedicated buffers for each client i : $\mathbf{in}_s(i)$ contains the rounds received from i , and $\mathbf{out}_s(i)$ the segments that have been sent to i .

A client is represented by a term $\langle P, \mathbf{s}, \delta_{\top}, \delta_{\mathbf{P}}, n, \mathbf{r}, \mathbf{in}_c \rangle_i$. As for GSP, i is its identity and P is its program. Note that the language for programs remains unaltered. The component δ_{\top} is analogous to \mathbf{u}_{\top} in the GSP model, i.e., it keeps all local updates until the client performs `push`. Differently, $\delta_{\mathbf{P}}$ keeps all finished blocks that have not been sent. The number n identifies the current round. Buffer \mathbf{r} keeps all sent rounds that have not been confirmed by the global store (similar to $\mathbf{b}_{\mathbf{P}}$ in GSP), while \mathbf{in}_c keeps all received segments (analogous to j).

We also impose the following well-formedness condition on systems.

Definition 3.2 (IGSP well-formedness). *A IGSP system $\mathbb{N} = \mathbf{C}_0 \parallel \dots \parallel \mathbf{C}_m \parallel \mathbf{S}$ with $\mathbf{S} = \langle \mathbf{s}, \mathbf{f}, \mathbf{in}_s, \mathbf{out}_s \rangle$ and $\mathbf{C}_l = \langle P_l, \mathbf{s}_l, \delta_{\top l}, \delta_{\mathbf{P}l}, n_l, \mathbf{r}_l, \mathbf{in}_{cl} \rangle_{i_l}$ for $l \in \{0, \dots, m\}$ is well-formed if the following conditions hold*

1. $i_l \neq i_{l'}$ for all $l \neq l'$.
2. $\text{dom}(\mathbf{in}_s) = \text{dom}(\mathbf{out}_s) \subseteq \{0, \dots, m\}$.
3. $i_l \notin \text{dom}(\mathbf{out}_s)$ implies $\mathbf{in}_{cl} = \epsilon$.
4. if $i_l \in \text{dom}(\mathbf{out}_s)$ and $\mathbf{in}_{cl} \cdot \mathbf{out}_s(i_l) \neq \epsilon$ then either
 - (i) $\mathbf{in}_{cl} \cdot \mathbf{out}_s(i_l) = \langle \delta_0, \mathbf{f}_0 \rangle \cdots \langle \delta_h, \mathbf{f}_h \rangle$; or
 - (ii) $\mathbf{in}_{cl} \cdot \mathbf{out}_s(i_l) = \langle \mathbf{s}, \mathbf{f}_0 \rangle \cdot \langle \delta_1, \mathbf{f}_1 \rangle \cdots \langle \delta_h, \mathbf{f}_h \rangle$
and $\mathbf{f}_j(i_l) \leq \mathbf{f}_k(i_l)$ for all $j < k \in \{0..h\}$ and $\mathbf{f}_h = \mathbf{f}$.
5. $\mathbf{r}_l = \langle i_0, n_0, \delta_0 \rangle \cdots \langle i_r, n_r, \delta_r \rangle$, $n_j < n_k$ for all $j < k \in \{0..r\}$ and either
 - (i) $\delta_{\mathbf{P}l} = \delta_0$, $n_r \leq n_l$ and $\mathbf{f}(i_l) \leq n_l$; or
 - (ii) $\delta_{\mathbf{P}l} \neq \delta_0$, $n_r < n_l$ and $\mathbf{f}(i_l) < n_l$.
6. either
 - (i) $\mathbf{r}_l = \epsilon$, $\mathbf{f}(i_l) \leq n_l$ and if $i_l \in \text{dom}(\mathbf{in}_s)$ then $\mathbf{in}_s(i_l) = \epsilon$;
 - (ii) $\mathbf{r}_l = \mathbf{in}_s(i_l) = \langle i_l, n_{fst}, \delta_{fst} \rangle \cdot \mathbf{r}'_l$ and $n_{fst} > \mathbf{f}(i_l)$;
 - (iii) $\mathbf{r}_l = \mathbf{r}''_l \cdot \langle i_l, n_{lst}, \delta_{lst} \rangle \cdot \mathbf{in}_s(i_l)$ and $\mathbf{in}_{cl} \cdot \mathbf{out}_s(i_l) = \langle \delta_0, \mathbf{f}_0 \rangle \cdots \langle \delta'_{lst}, \mathbf{f}_{lst} \rangle$ with $\mathbf{f}_{lst}(i_l) = n_{lst}$; or
 - (iv) $\mathbf{in}_s(i_l) = \epsilon$, $\mathbf{r}_l = \mathbf{r}''_l \cdot \langle i_l, n_{lst}, \delta_{lst} \rangle$ and either $\mathbf{in}_{cl} \cdot \mathbf{out}_s(i_l) = \langle \mathbf{s}, \mathbf{f}' \rangle \cdot \mathbf{seg}$ with $\mathbf{f}'(i_l) \leq n_{lst}$ or $\mathbf{in}_{cl} \cdot \mathbf{out}_s(i_l) = \epsilon$ and $\mathbf{f}(i_l) \leq n_{lst}$.

We require all clients to have different identifiers (1). Communication channels in the implementation are bidirectional, hence $i_l \in \text{dom}(\mathbf{in}_s)$ iff $i_l \in \text{dom}(\mathbf{out}_s)$ (2). Moreover, the input buffer of a disconnected client is empty (3). Condition (4) states that $\langle \mathbf{s}, \mathbf{f}' \rangle$ can appear only as the first message in the flow from the store and that the store confirms processed rounds in a non-decreasing order. Similarly, clients send rounds with increasing round number (5). The last condition (6) states a coherence requirement between pending rounds and the segments sent by the store, which can only confirm rounds that are pending.

$$\begin{array}{c}
\text{(I-UPDATE)} \\
\langle \text{update}(u); P, \mathbf{s}, \delta_{\top}, \delta_{\text{P}}, n, \mathbf{r}, \text{inc} \rangle_i \parallel \mathbf{N} \xrightarrow{wr(u^V)}_i \langle P, \mathbf{s}, \text{append}(\delta_{\top}, u), \delta_{\text{P}}, n, \mathbf{r}, \text{inc} \rangle_i \parallel \mathbf{N} \\
\\
\text{(I-PUSH)} \\
\langle \text{push}; P, \mathbf{s}, \delta_{\top}, \delta_{\text{P}}, n, \mathbf{r}, \text{inc} \rangle_i \parallel \mathbf{N} \xrightarrow{push}_i \langle P, \mathbf{s}, \delta_{\emptyset}, \text{reduce}(\delta_{\text{P}} \cdot \delta_{\top}), n+1, \mathbf{r}, \text{inc} \rangle_i \parallel \mathbf{N} \\
\\
\text{(I-SEND)} \\
\frac{\delta_{\text{P}} \neq \delta_{\emptyset} \quad i \in \text{dom}(\text{ins}) \quad r = \langle i, n, \delta_{\text{P}} \rangle \quad \text{inc} \cdot \text{outs}(i) \neq \langle \mathbf{s}', \mathbf{f} \rangle \cdot \text{seg}}{\langle P, \mathbf{s}, \delta_{\top}, \delta_{\text{P}}, n, \mathbf{r}, \text{inc} \rangle_i \parallel \langle \mathbf{s}', \mathbf{f}, \text{ins}, \text{outs} \rangle_i \parallel \mathbf{C} \xrightarrow{\tau}_i \langle P, \mathbf{s}, \delta_{\top}, \delta_{\emptyset}, n, \mathbf{r} \cdot r, \text{inc} \rangle_i \parallel \langle \mathbf{s}', \mathbf{f}, \text{ins}[i \mapsto \text{ins}(i) \cdot r], \text{outs} \rangle_i \parallel \mathbf{C}} \\
\\
\text{(I-RECEIVE)} \\
\frac{\text{outs}(i) = \text{seg} \cdot \text{seg}}{\langle P, \mathbf{s}, \delta_{\top}, \delta_{\text{P}}, n, \mathbf{r}, \text{inc} \rangle_i \parallel \langle \mathbf{s}', \mathbf{f}, \text{ins}, \text{outs} \rangle_i \parallel \mathbf{C} \xrightarrow{\tau}_i \langle P, \mathbf{s}, \delta_{\top}, \delta_{\text{P}}, n, \mathbf{r}, \text{inc} \cdot \text{seg} \rangle_i \parallel \langle \mathbf{s}', \mathbf{f}, \text{ins}, \text{outs}[i \mapsto \text{seg}] \rangle_i \parallel \mathbf{C}} \\
\\
\text{(I-PULL}_1\text{)} \\
\frac{\mathbf{r}' = \text{filter}(\mathbf{f}_k(i), \mathbf{r}) \quad \text{inc} = \langle \delta_1, \mathbf{f}_1 \rangle \dots \langle \delta_k, \mathbf{f}_k \rangle}{\langle \text{pull}; P, \mathbf{s}, \delta_{\top}, \delta_{\text{P}}, n, \mathbf{r}, \text{inc} \rangle_i \parallel \mathbf{N} \xrightarrow{pull}_i \langle P, \text{apply}(\mathbf{s}, \text{reduce}(\delta_1 \dots \delta_k)), \delta_{\top}, \delta_{\text{P}}, n, \mathbf{r}', \epsilon \rangle_i \parallel \mathbf{N}} \\
\\
\text{(I-READ)} \\
\frac{\text{read}(r, \text{apply}(\mathbf{s}, \Delta(\mathbf{r}) \cdot \delta_{\text{P}} \cdot \delta_{\top})) = v}{\langle \text{let } x = \text{read}(r) \text{ in } P, \mathbf{s}, \delta_{\top}, \delta_{\text{P}}, n, \mathbf{r}, \text{inc} \rangle_i \parallel \mathbf{N} \xrightarrow{rd(r)}_i \langle P\{v/x\}, \mathbf{s}, \delta_{\top}, \delta_{\text{P}}, n, \mathbf{r}, \text{inc} \rangle_i \parallel \mathbf{N}} \\
\\
\text{(I-CONFIRM)} \\
\frac{v = (\mathbf{r} \cdot \delta_{\text{P}} \cdot \delta_{\top}) == \epsilon}{\langle \text{let } x = \text{confirmed in } P, \mathbf{s}, \delta_{\top}, \delta_{\text{P}}, n, \mathbf{r}, \text{inc} \rangle_i \parallel \mathbf{N} \xrightarrow{cfm}_i \langle P\{v/x\}, \mathbf{s}, \delta_{\top}, \delta_{\text{P}}, n, \mathbf{r}, \text{inc} \rangle_i \parallel \mathbf{N}} \\
\\
\text{(I-BATCH)} \\
\frac{\langle \delta, \mathbf{f}' \rangle = \text{rmds}(\text{ins}) \quad \delta \neq \delta_{\emptyset} \quad \mathbf{s}' = \text{apply}(\mathbf{s}, \delta) \quad \forall i. (\text{outs}'(i) = \text{outs}(i) \cdot \langle \delta, \mathbf{f}[i] \rangle) \wedge \text{ins}'(i) = \epsilon}{\langle \mathbf{s}, \mathbf{f}, \text{ins}, \text{outs} \rangle_i \parallel \mathbf{C} \xrightarrow{\tau} \langle \mathbf{s}', \mathbf{f}[\mathbf{f}'], \text{ins}', \text{outs}' \rangle_i \parallel \mathbf{C}} \\
\\
\text{(I-DROP-CXN)} \\
\frac{i \in \text{ins} \quad i \in \text{outs}}{\langle P, \mathbf{s}, \delta_{\top}, \delta_{\text{P}}, n, \mathbf{r}, \text{inc} \rangle_i \parallel \langle \mathbf{s}, \mathbf{f}, \text{ins}, \text{outs} \rangle_i \parallel \mathbf{C} \xrightarrow{\tau} \langle P, \mathbf{s}, \delta_{\top}, \delta_{\text{P}}, n, \mathbf{r}, \epsilon \rangle_i \parallel \langle \mathbf{s}, \mathbf{f}, \text{ins} \setminus i, \text{outs} \setminus i \rangle_i \parallel \mathbf{C}} \\
\\
\text{(I-ACCEPT-CXN)} \\
\frac{i \notin \text{ins} \quad i \notin \text{outs}}{\langle \mathbf{s}, \mathbf{f}, \text{ins}, \text{outs} \rangle_i \parallel \mathbf{C}_i \parallel \mathbf{C} \xrightarrow{\tau} \langle \mathbf{s}, \mathbf{f}, \text{ins}[i \mapsto \epsilon], \text{outs}[i \mapsto \langle \mathbf{s}, \mathbf{f} \rangle] \rangle_i \parallel \mathbf{C}_i \parallel \mathbf{C}} \\
\\
\text{(I-PULL}_2\text{)} \\
\frac{\text{inc} = \langle \mathbf{s}''', \mathbf{f}_0 \rangle \cdot \langle \delta_1, \mathbf{f}_1 \rangle \dots \langle \delta_k, \mathbf{f}_k \rangle \quad \mathbf{s}'' = \text{apply}(\mathbf{s}''', \text{reduce}(\delta_{\emptyset} \cdot \delta_1 \dots \delta_k)) \quad \mathbf{r}' = \text{filter}(\mathbf{f}_k(i), \mathbf{r})}{\langle \text{pull}; P, \mathbf{s}, \delta_{\top}, \delta_{\text{P}}, n, \mathbf{r}, \text{inc} \rangle_i \parallel \langle \mathbf{s}', \mathbf{f}, \text{ins}, \text{outs} \rangle_i \parallel \mathbf{C} \xrightarrow{pull}_i \langle P, \mathbf{s}'', \delta_{\top}, \delta_{\text{P}}, n, \mathbf{r}', \epsilon \rangle_i \parallel \langle \mathbf{s}', \mathbf{f}, \text{ins}[i \mapsto \text{ins}(i) \cdot \mathbf{r}'], \text{outs} \rangle_i \parallel \mathbf{C}}
\end{array}$$

Fig. 4. Operational semantics of IGSP

3.2 Operational Semantics

As for the idealised model, the operational semantics is given by a labelled transition system over well-formed terms, up-to structural equivalence. We consider a new label τ without any client annotation for transitions associated with changes in the global store and communication failures. The inference rules are in Fig. 4.

Rule (I-UPDATE), which is analogous to rule (UPDATE), adds the operation u to the temporary block δ_τ . The decoration \mathbf{v} is irrelevant in this model, hence we do not impose any freshness requirement. A client terminates a block by executing `push` (I-PUSH). At this time, the block δ_τ is appended to the already terminated blocks in δ_P , which will be sent on the next round. Additionally, the block counter n is incremented by 1. By rule (I-SEND), a client sends changes to the global store. This transition takes place whenever the client is connected (i.e., $i \in \text{dom}(\mathbf{in}_s)$), there are finished blocks in δ_P (i.e., $\delta_P \neq \delta_\emptyset$) and there is no need for resynchronisation (i.e., $\mathbf{in}_c \cdot \text{out}_s(i) \neq \langle \mathbf{s}', \mathbf{f} \rangle \cdot \mathbf{seg}$)³. The available blocks are sent within the same round $r = \langle i, n, \delta_P \rangle$, which contains the number n corresponding to the last finished block. The new round r is added to the corresponding input buffer in the store, i.e., $\mathbf{in}_s(i)$ is updated to $\mathbf{in}_s(i) \cdot r$ (where $[- \mapsto -]$ is the update operator for functions). Additionally, r is added to the sequence of pending rounds \mathbf{r} and the buffer δ_P is reset to δ_\emptyset .

Symmetrically, the client i may receive an available segment at any time (I-RECEIVE). The new segment seg is removed from the buffer $\text{out}_s(i)$ of the global store and added to the input buffer of the client. As for the idealised model, all received changes are applied to the local replica when i performs `pull`. Rule (I-PULL₁) handles the case in which the connection with the global store has not been previously reset. In such case, all received segments are of the form $\langle \delta, \mathbf{f} \rangle$. Therefore, the changes $\delta_1 \cdots \delta_k$ are applied to the local state \mathbf{s} and all rounds confirmed by the received segments are removed from the pending list \mathbf{r} . By well-formedness (Def. 3.2,5), it suffices to consider the confirmation \mathbf{f}_k , which has the greatest confirmation. Hence, all rounds up-to $\mathbf{f}_k(i)$ are removed from \mathbf{r} . This is done by the auxiliary function $filter(-, -)$, defined as follows

$$filter(n, \mathbf{r}) = \langle i, n_j, \delta_j \rangle \cdots \langle i, n_k, \delta_k \rangle \quad \text{if } \mathbf{r} = \langle i, n_0, \delta_0 \rangle \cdots \langle i, n_j, \delta_j \rangle \cdots \langle i, n_k, \delta_k \rangle, \\ n_{j-1} \leq n \text{ and } n_j > n$$

Rules (I-READ) and (I-CONFIRM) are analogous the ones in the GSP calculus. We use $\Delta(-)$ for the function that projects a sequence of rounds into the sequence that contains the corresponding δ s. The global store changes its state as prescribed by rule (I-BATCH): it collects all received rounds in \mathbf{in}_s by using the auxiliary function $rnds(-)$, which builds a unique object δ by appending all available rounds, and a function \mathbf{f} that associates each client with the number of the last received round. Let \mathbf{in}_s be such that $\text{dom}(\mathbf{in}_s) = \{i_0, \dots, i_m\}$ and $\forall i_l \in \text{dom}(\mathbf{in}_s). \mathbf{in}_s(i_l) = \mathbf{r}_l \cdot \langle i_l, n_l^{k_l}, \delta_l^{k_l} \rangle$. Then, $rnds(-)$ is defined as follows

³ For simplicity we check re-synchronisation by inspecting buffers instead of explicitly adding the condition *channel established* used in the implementation.

$$\begin{aligned} rns(\mathbf{in}_s) = \langle \delta, \mathbf{f}' \rangle \quad \text{with} \quad & \delta = \text{reduce}(\Delta(\mathbf{in}_s(i_0)) \cdots \Delta(\mathbf{in}_s(i_m))), \\ & \text{dom}(\mathbf{f}') = \{i \mid i \in \text{dom}(\mathbf{in}_s) \text{ and } \mathbf{in}_s(i) \neq \epsilon\} \text{ and} \\ & \forall i_l \in \text{dom}(\mathbf{in}_s). \mathbf{f}'(i_l) = n_l^{k_l} \end{aligned}$$

The obtained δ is applied to the current state \mathbf{s} and \mathbf{f} is updated with \mathbf{f}' . In addition, the new segment $\langle \delta, \mathbf{f}' \rangle$ is sent to every connected client, i.e., it is added at the end of every buffer $\text{out}_s(i)$. The input buffers $\mathbf{in}_s(i_l)$ are emptied because all received rounds have been processed.

The remaining rules deal with connectivity issues: rule (I-DROP-CXN) models a disconnection: the buffers $\text{out}_s(i)$ and $\mathbf{in}_s(i)$ are removed from the global store and also the input buffer of i is set to ϵ . When the client i (re-)establishes its connection (I-ACCEPT-CXN), the store creates the buffers for i and sends a segment containing the current state of the store. Rule (I-PULL₂) is analogous to (I-PULL₁), but handles the first pull after a reconnection. The first received segment $\langle \mathbf{s}''', \mathbf{f}_0 \rangle$ contains a state instead of a delta object. The client uses \mathbf{s}''' instead of its local state to resynchronise. The application of successive segments is analogous to rule (I-PULL₁). Moreover, the client resends a round \mathbf{r}' containing all pending segments lost by the server during the disconnection.

The proposed implementation allows for a server to crash, i.e., to close all communication buffers, but we do not model explicitly this behaviour because it can be obtained by applying rule (I-DROP-CXN) several times.

Lemma 3.1. *Let N be a well-formed IGSP system. If $N \xrightarrow{\lambda}_i N'$, then N' is well-formed.*

4 Correctness of the Implementation

We now prove that IGSP is a correct implementation of GSP. We recall in Fig. 5 the requirements stated in [5] for the operations provided by the data types State and Delta. Formally, the relation \triangleleft associates delta and state objects with sequences of updates: $\delta \triangleleft \mathbf{u}$ (similarly, $\mathbf{s} \triangleleft \mathbf{u}$) means that δ (correspondingly, \mathbf{s}) is a compact representation of \mathbf{u} . Then, it is also assumed that $\mathbf{s} \triangleleft \mathbf{u}$ implies $\text{read}(r, \mathbf{s}) = \text{rvalue}(r, \mathbf{u})$ for any r . Building on the above relation, we define under which conditions a IGSP system is an implementation of a GSP system.

Definition 4.1. *Let $N = C_0 \parallel \dots \parallel C_m \parallel S$ be a GSP system such that $C_l = \langle P_l, \mathbf{u}_{\top l}, \mathbf{b}_{S_l}, \mathbf{b}_{P_l}, k_l, j_l \rangle_{i_l}$ for all $l \in \{0, \dots, m\}$, and $N = \mathbf{C}_0 \parallel \dots \parallel \mathbf{C}_m \parallel S$ a IGSP system such that $\mathbf{S} = \langle \mathbf{s}, \mathbf{f}_s, \mathbf{in}_s, \text{out}_s \rangle$ and $\mathbf{C}_l = \langle P_l, \mathbf{s}_l, \delta_{\top l}, \delta_{P_l}, n_l, \mathbf{r}_l, \mathbf{in}_{cl} \rangle_{i_l}$. We say N implements N if the following conditions hold:*

1. $\mathbf{s} \triangleleft \underline{S}$;
2. $\mathbf{s}_l \triangleleft \underline{S[0..k_l - 1]}$;
3. $\delta_{\top l} \triangleleft \underline{\mathbf{u}_{\top l}}$;
4. $\text{reduce}(\Delta(\mathbf{r}_l) \cdot \delta_{P_l}) \triangleleft \underline{\mathbf{b}_{P_l}}$;

$$\begin{array}{c}
 \begin{array}{ccc}
 (\leftarrow \delta_\emptyset) & \begin{array}{c} (\leftarrow \text{append}) \\ \delta \triangleleft \mathbf{u} \end{array} & \begin{array}{c} (\leftarrow \text{read}) \\ \mathbf{s} \triangleleft \mathbf{u} \end{array} \\
 \delta_\emptyset \triangleleft \epsilon & \frac{}{\text{append}(\delta, \mathbf{u}) \triangleleft \mathbf{u} \cdot \mathbf{u}'} & \frac{}{\text{read}(r, \mathbf{s}) = \text{rvalue}(r, \mathbf{u})} \\
 & (\leftarrow \text{apply}) & (\leftarrow \text{reduce}) \\
 (\leftarrow \emptyset) & \frac{\mathbf{s} \triangleleft \mathbf{u} \quad \delta_1 \triangleleft \mathbf{u}_1 \dots \delta_n \triangleleft \mathbf{u}_n}{\text{apply}(\mathbf{s}, \delta_1 \dots \delta_n) \triangleleft \mathbf{u} \cdot \mathbf{u}_1 \dots \mathbf{u}_n} & \frac{\delta_1 \triangleleft \mathbf{u}_1 \dots \delta_n \triangleleft \mathbf{u}_n}{\text{reduce}(\delta_1 \dots \delta_n) \triangleleft \mathbf{u}_1 \dots \mathbf{u}_n} \\
 \emptyset \triangleleft \epsilon & &
 \end{array}
 \end{array}$$

Fig. 5. Coherence requirements for *Delta* and *State* operators

5. if $i_l \in \text{dom}(\text{in}_s)$ and $\text{in}_c \cdot \text{out}_s(i_l) \neq \langle \mathbf{s}', \mathbf{f} \rangle \cdot \text{seg}$ then
 $\text{reduce}(\Delta(\text{in}_s(i_l)) \cdot \delta_{\mathbf{P}_l}) \triangleleft \mathbf{b}_{\mathbf{S}_l}$;
6. if $i_l \in \text{dom}(\text{out}_s)$ then either
 - i. $\text{in}_{cl} \cdot \text{out}_s(i_l) = \epsilon$, $k_l = |S|$;
 - ii. $\text{in}_{cl} = \langle \delta', \mathbf{f} \rangle \cdot \text{seg}$, $\text{reduce}(\Delta(\text{in}_{cl})) \triangleleft S[k_l..k_l + j_l - 1]$, and
 $\text{reduce}(\Delta(\text{out}_s(i_l))) \triangleleft S[k_l + j_l..|S| - 1]$;
 - iii. $\text{in}_{cl} = \langle \mathbf{s}', \mathbf{f}_0 \rangle \cdot \text{seg}$, there exists t s.t. $k_l \leq t \leq k_l + j_l$ s.t. $\mathbf{s}' \triangleleft S[0..t - 1]$,
 $\text{reduce}(\Delta(\text{seg})) \triangleleft S[t..k_l + j_l - 1]$ and
 $\text{reduce}(\Delta(\text{out}_s(i_l))) \triangleleft S[k_l + j_l..|S| - 1]$; or
 - iv. $\text{in}_{cl} = \epsilon$, $\text{out}_s(i_l) = \langle \mathbf{s}', \mathbf{f} \rangle \cdot \text{seg}$ there exists $t \geq k_l + j_l$ s.t. $\mathbf{s}' \triangleleft S[0..t - 1]$
 $\text{reduce}(\Delta(\text{seg})) \triangleleft S[t..|S| - 1]$;
7. for all \mathbf{f} s.t. $\mathbf{f} = \mathbf{f}_s$ or $\langle \cdot, \mathbf{f} \rangle \in \text{in}_{cl} \cdot \text{out}_s(i_l)$, for all $\langle i, n, \delta' \rangle \in \mathbf{r}_l$ if $n \leq \mathbf{f}(i_l)$
then $\delta \triangleleft S[x..x']$ and $\delta' \triangleleft S[y..y']$ with $y' \leq x'$.

The first three conditions are self-explanatory. Condition (4) states that the pending blocks in $\mathbf{b}_{\mathbf{P}_l}$ correspond either to rounds in the pending list \mathbf{r}_l or to blocks ready to be sent, i.e., in $\delta_{\mathbf{P}_l}$. By condition (5), if a client is synchronised with the store (i.e., $i_l \in \text{dom}(\text{in}_s)$ and $\text{in}_c \cdot \text{out}_s(i_l) \neq \langle \mathbf{s}', \mathbf{f} \rangle \cdot \text{seg}$) then all blocks in the sending list $\mathbf{b}_{\mathbf{S}_l}$ are either rounds that have been sent, i.e., in $\text{in}_s(i_l)$, or ready blocks in $\delta_{\mathbf{P}_l}$. Condition (6) establishes the relation between the received messages in both models. Basically, the local replica is complete when there are no segments for the client (i). When the first received segment is a delta object (ii), the content in the input buffer in_{cl} corresponds to the received messages in $S[k_l..k_l + j_l - 1]$ and the the output buffer $\text{out}_s(i_l)$ contains the updates in the sequence $S[k_l + j_l..|S| - 1]$. In the remaining two cases, the first segment contains a state. When the segment is in the input buffer of the client (iii), the received state \mathbf{s}' corresponds to a prefix of the sequence S whose length lies in between of the updates already received by the client in the idealised model, i.e., $S[0..t - 1]$ with $k \leq t \leq k_l + j_l$, while the remaining conditions are analogous to the previous case. Differently, when the first segment is still on the output buffer of the store (iv), \mathbf{s}' corresponds to a prefix that contains at least all updates in S already known to the client, i.e., $t \geq k_l + j_l$ because the store confirmations are monotonic.

Condition (7) states that in any segment $\langle \delta, \mathbf{f} \rangle$ sent by the store, δ corresponds to a contiguous sequence of updates in S , i.e., $S[x..x']$. Moreover, all confirmed rounds are also within the prefix $S[0..x']$.

We now show that IGSP is a correct implementation of GSP by proving that N weakly simulates \mathbb{N} when \mathbb{N} implements N . We use standard simulation but technically we take into account the fact that GSP associates a fresh event identifier to each update while IGSP does not. Take $\rightarrow = \xrightarrow{\tau} \bigcup_{i \in \mathcal{I}} \xrightarrow{\tau}_i, \Rightarrow$ as the reflexive and transitive closure of \rightarrow , i.e. $\Rightarrow = \rightarrow^*$, and $\xRightarrow{\lambda}_i = \Rightarrow; \xrightarrow{\lambda}_i; \Rightarrow$.

Definition 4.2 (Simulation). \mathcal{R} is an implementation simulation if for all $(\mathbb{N}, N) \in \mathcal{R}$ we have:

1. If $\mathbb{N} \xrightarrow{wr(u^v)}_i N'$ then $\exists N', w$ s.t. $N \xrightarrow{wr(u^w)}_i N'$ and $(N', N') \in \mathcal{R}$;
2. If $\mathbb{N} \xrightarrow{\lambda}_i N'$ and $\lambda \neq wr(u^v), \tau$ then $\exists N'$ s.t. $N \xRightarrow{\lambda}_i N'$ and $(N', N') \in \mathcal{R}$;
3. If $\mathbb{N} \rightarrow N'$ then $\exists N'$ s.t. $N \Rightarrow N'$ and $(N', N') \in \mathcal{R}$;

As usual, we write $\mathbb{N} \lesssim N$ if there exists a simulation \mathcal{R} s.t. $(\mathbb{N}, N) \in \mathcal{R}$.

Theorem 4.1. If \mathbb{N} implements N , then $\mathbb{N} \lesssim N$.

Proof. We show that $\mathcal{R} = \{(\mathbb{N}, N) \mid \mathbb{N} \text{ implements } N\}$ is a simulation.

We remark that \mathcal{R}^{-1} is not a simulation because the implementation cannot mimic the behaviour in which a client have completed two consecutive blocks (i.e., two `push` commands) without sending the first block. In GSP it is still possible to interleave the two blocks with blocks sent by other clients but in IGSP they are treated as atomic because they will be sent as a unique δ object.

5 Consistency Guarantees

In this section we study the consistency properties offered by GSP. We rely on the characterisation of properties in terms of abstract executions [4], execution histories enriched with information about visibility and arbitration of actions.

Definition 5.1. Let N be a well-formed GSP system, an abstract history for N is a tuple $\mathbb{A} = \langle N, \text{OP}, \text{SS}, \text{SO}, \text{VIS}, \text{AR} \rangle$ where:

- $\text{OP} : \mathbb{V} \rightarrow \mathcal{R} \cup \mathcal{U}$ maps events to operations;
- $\text{SS} : \mathcal{I} \rightarrow \mathbb{V}$ associates events with a session (i.e., a client);
- $\text{SO} \subseteq \mathbb{V} \times \mathbb{V}$ describes the order of operations within a session;
- $\text{VIS} \subseteq \mathbb{V} \times \mathbb{V}$ indicates whether the effects of an update are visible to a read;
- $\text{AR} \subseteq \mathbb{V} \times \mathbb{V}$ resolves concurrent update conflicts.

We write \downarrow_{\cdot} for function/relation restriction. For a given abstract history \mathbb{A} , we write \mathbb{U} (similarly, \mathbb{R}) for the codomain restriction of OP to \mathcal{U} (correspondingly, \mathcal{R}), i.e., $\mathbb{U} = \{\mathbf{v} \mid \mathbf{v} \in \text{OP}, \text{OP}(\mathbf{v}) \in \mathcal{U}\}$ ($\mathbb{R} = \{\mathbf{v} \mid \mathbf{v} \in \text{OP}, \text{OP}(\mathbf{v}) \in \mathcal{R}\}$).

Definition 5.2 (Well-formed history). Let $N = C_0 \parallel \dots \parallel C_m \parallel S$ be a GSP system where $C_l = \langle P_l, \mathbf{u}_{Tl}, \mathbf{b}_{Sl}, \mathbf{b}_{Pl}, k_l, j_l \rangle_i$. A history $\mathbb{A} = \langle N, \text{OP}, \text{SS}, \text{SO}, \text{VIS}, \text{AR} \rangle$ is well-formed if the following conditions hold:

1. for all $i \in \text{dom}(\text{ss})$, $\text{ss}(i) \subseteq \text{dom}(\text{op})$;
2. $(\mathbf{v}, \mathbf{w}) \in \text{so}$ then exist $i \in \text{dom}(\text{ss})$ s.t. $\{\mathbf{v}, \mathbf{w}\} \in \text{ss}(i)$.
3. for all $i \in \text{dom}(\text{ss})$, $\text{so} \downarrow_{\text{SS}(i)}$ is a total order;
4. $\text{vis} \subseteq \mathbb{U} \times \mathbb{R}$;
5. $\text{AR} \subseteq \mathbb{U} \times \mathbb{U}$ is a prefix order.
6. $(\mathbf{v}, \mathbf{w}) \in \text{AR}$ iff
 - $\underline{S}[i] = u^{\mathbf{v}}$ and $\underline{S}[j] = u^{\mathbf{w}}$ and $i < j$; or
 - $u^{\mathbf{v}} \in \underline{S}$ and $u^{\mathbf{w}} \in \underline{\mathbf{b}_{S_I} \cdot \mathbf{u}_{T_I}}$;
7. if $\underline{\mathbf{b}_{S_I} \cdot \mathbf{u}_{T_I}}[i] = u^{\mathbf{v}}$ and $\underline{\mathbf{b}_{S_I} \cdot \mathbf{u}_{T_I}}[j] = u^{\mathbf{w}}$ and $i < j$, then $(\mathbf{v}, \mathbf{w}) \in \text{so}$ and $\{\mathbf{v}, \mathbf{w}\} \in \text{ss}(i)$.

The above conditions ensure that events in ss are associated with an operation by op (1). Besides, so only relates events belonging to the same session (2), which are totally ordered within each session (3). Differently from the definition in [2], we restrict visibility to keep track of dependencies between updates and read events(4). We do not require AR to be a total order but instead to be a prefix order (5). In this way the updates in different replicas are arbitrated when they reach the global store. The remaining two conditions require the abstract history to be consistent with the state of the system.

Rules in Fig. 6 provides an operational way to associate abstract executions with GSP computations. Rules (A-UPDATE) and (A-READ) add new events to the history and corresponds to the execution of a read or update operation by a client. In both cases op is extended with a new event \mathbf{v} (i.e., $\mathbf{v} \notin \text{dom}(\text{op})$), which is associated with the corresponding operation (either r or u). The new event \mathbf{v} is added to the corresponding session i , and so is updated to make \mathbf{v} the maximal event for the session i . Rule (A-UPDATE) amends AR by capturing the fact that all updates that are already in the global state took place before the new event. Rule (A-READ) instead augments vis with the pairs associating the new event with all events that are seen by the read action, namely, the local view of the global state $S[0..k_i - 1]$ and the local buffers \mathbf{b}_{p_i} and \mathbf{u}_{T_i} . Rule (A-ARB) handles the changes in the state of the global store (due to a send transition in one client) and amends AR by arbitrating (i) the new events by respecting the relative order in which they are added to the store (i.e., $\{(\mathbf{v}_i, \mathbf{v}_j) \mid i, j \in \{0, \dots, n\}, i < j\}$) and (ii) all updates in the local state of the clients after the new ones (i.e., $(\{\mathbf{v}_0, \dots, \mathbf{v}_n\} \times \{\mathbf{w} \mid \mathbf{w} \in \mathbb{U} \forall u.u^{\mathbf{w}} \notin S \cdot \mathbf{u}\})$). The remaining transitions of the system are considered as internal changes that do not affect the history and are handled by rule (A-INT).

Lemma 5.1. *Let \mathbb{A} be a well-formed history. If $\mathbb{A} \xrightarrow{\lambda}_i \mathbb{A}'$, then \mathbb{A}' is well-formed.*

We use histories to analyse the ordering guarantees offered by the GSP model. (Due to space limitation, we refer the interested reader to see the characterisations provided in [2, Ch. 5]).

Theorem 5.1. *If $\langle N, \emptyset, \emptyset, \emptyset, \emptyset, \emptyset \rangle \rightarrow_i^* \langle N', \text{op}, \text{ss}, \text{so}, \text{vis}, \text{AR} \rangle$ then*

- (1) *Read My Writes: $\text{so} \downarrow_{\mathbb{U} \times \mathbb{R}} \subseteq \text{vis}$*

$$\begin{array}{c}
\text{(A-UPDATE)} \\
\frac{N \xrightarrow{wr(u^v)}_i N' \quad v \notin \text{dom}(\text{OP}) \quad \text{OP}' = \text{OP}[v \mapsto u] \quad \text{SS}' = \text{SS}[i \mapsto \text{SS}(i) \cup \{v\}]}{\text{SO}' = \text{SO} \cup (\text{SS}(i) \times \{v\}) \quad \text{AR}' = \text{AR} \cup (\{\mathbf{w} \mid u^{\mathbf{w}} \in S\} \times \{v\})} \\
\langle N, \text{OP}, \text{SS}, \text{SO}, \text{VIS}, \text{AR} \rangle \xrightarrow{wr(u^v)}_i \langle N', \text{OP}', \text{SS}, \text{SO}', \text{VIS}, \text{AR}' \rangle \\
\text{(A-READ)} \\
\frac{N \xrightarrow{rd(r)}_i N' \quad v \notin \text{dom}(\text{OP}) \quad \text{OP}' = \text{OP}[v \mapsto r] \quad \text{SS}' = \text{SS}[i \mapsto \text{SS}(i) \cup \{v\}]}{\text{SO}' = \text{SO} \cup (\text{SS}(i) \times \{v\}) \quad \text{VIS}' = \text{VIS} \cup (\{\mathbf{w} \mid u^{\mathbf{w}} \in S[0..k_i - 1] \cdot \mathbf{b}_{P_i} \cdot \mathbf{u}_{T_i}\} \times \{v\})} \\
\langle N, \text{OP}, \text{SS}, \text{SO}, \text{VIS}, \text{AR} \rangle \xrightarrow{rd(r^v)}_i \langle N', \text{OP}', \text{SS}', \text{SO}', \text{VIS}', \text{AR}' \rangle \\
\text{(A-ARB)} \\
\frac{S \parallel C \xrightarrow{\lambda}_i S \cdot \mathbf{u} \parallel C' \quad \lambda \neq wr(u^v), rd(r) \quad \mathbf{u} = (u_0^{v_0} \dots u_n^{v_n})}{\text{AR}' = \text{AR} \cup \{(\mathbf{v}_i, \mathbf{v}_j) \mid i, j \in \{0, \dots, n\}, i < j\} \cup (\{\mathbf{v}_0, \dots, \mathbf{v}_n\} \times \{\mathbf{w} \mid \mathbf{w} \in \mathbb{U}, \forall u.u^{\mathbf{w}} \notin S \cdot \mathbf{u}\})} \\
\langle S \parallel C, \text{OP}, \text{SS}, \text{SO}, \text{VIS}, \text{AR} \rangle \xrightarrow{\lambda}_i \langle S \cdot \mathbf{u} \parallel C', \text{OP}, \text{SS}, \text{SO}, \text{VIS}, \text{AR}' \rangle \\
\text{(A-INT)} \\
\frac{S \parallel C \xrightarrow{\lambda}_i S \parallel C' \quad \lambda \neq wr(u^v), rd(r)}{\langle S \parallel C, \text{OP}, \text{SS}, \text{SO}, \text{VIS}, \text{AR} \rangle \xrightarrow{\lambda}_i \langle S \parallel C', \text{OP}, \text{SS}, \text{SO}, \text{VIS}, \text{AR} \rangle}
\end{array}$$

Fig. 6. Computation of abstract executions

- (2) *Monotonic Read*: $\text{VIS}; \text{SO} \downarrow_{\mathbb{U} \times \mathbb{R}} \subseteq \text{VIS}$.
- (3) *No Circular Causality*: $(\text{SO} \cup \text{VIS})^+$ is acyclic.
- (4) *Causal Visibility*: $(\text{SO} \cup \text{VIS})^+ \downarrow_{\mathbb{U} \times \mathbb{R}} \subseteq \text{VIS}$.
- (5) *Causal Arbitration*: $((\text{SO} \cup \text{VIS})^+ \setminus \text{SO}) \downarrow_{\mathbb{U} \times \mathbb{U}} \subseteq \text{AR}$.
- (6) *Consistent prefix*: $\text{AR}; (\text{VIS} \setminus \text{SS}) \subseteq \text{VIS}$.

The following example shows that the GSP model exhibits the Dekker anomaly, hence it does not enjoy *sequential consistency* [2].

Example 5.1 (Dekker anomaly). Consider the following system consisting of two clients and the empty store $N = \epsilon \parallel C_1 \parallel C_2$ where

$$\begin{aligned}
C_1 &= \langle \text{update}(u_1); \text{let } y = \text{read}(r_1) \text{ in } P, \epsilon, \epsilon, \epsilon, 0, 0 \rangle_{i_1} \\
C_2 &= \langle \text{update}(u_2); \text{let } y = \text{read}(r_2) \text{ in } Q, \epsilon, \epsilon, \epsilon, 0, 0 \rangle_{i_2}
\end{aligned}$$

Since the updates are made locally, none of the clients see the update performed by the other and this is the essence of the Dekker anomaly which is ruled out by strong consistency models like sequential consistency or linearizability.

6 GSP with atomic updates

In this section we study the atomic updates proposed in [5]. We extend the language of programs as follows:

$$\text{(PROGRAM)} \quad P ::= \dots \mid \text{syncUpd}(u); P$$

$$\begin{array}{c}
\text{(SYNC-UPD)} \\
\langle \text{syncUpd}(u); P, u_{\top}, b_{\text{S}}, b_{\text{P}}, k, j \rangle_i \parallel N \xrightarrow{\tau}_i \langle \text{update}(u); \text{push}; \text{wait}; P, u_{\top}, b_{\text{S}}, b_{\text{P}}, k, j \rangle_i \parallel N \\
\text{(WAIT)} \\
\langle \text{wait}; P, u_{\top}, b_{\text{S}}, b_{\text{P}}, k, j \rangle_i \parallel N \xrightarrow{\tau}_i \langle \text{let } x = \text{confirmed in } x \triangleright (\text{pull}; \text{wait}; P); P, u_{\top}, b_{\text{S}}, b_{\text{P}}, k, j \rangle_i \parallel N \\
\text{(GUARD-TRUE)} \\
\frac{e \downarrow \text{true}}{\langle e \triangleright (P); Q, u_{\top}, b_{\text{S}}, b_{\text{P}}, k, j \rangle_i \parallel N \xrightarrow{\tau}_i \langle Q, u_{\top}, b_{\text{S}}, b_{\text{P}}, k, j \rangle_i \parallel N} \\
\text{(GUARD-FALSE)} \\
\frac{e \downarrow \text{false}}{\langle e \triangleright (P); Q, u_{\top}, b_{\text{S}}, b_{\text{P}}, k, j \rangle_i \parallel N \xrightarrow{\tau}_i \langle P, u_{\top}, b_{\text{S}}, b_{\text{P}}, k, j \rangle_i \parallel N}
\end{array}$$

Fig. 7. Semantics of GSP with atomic updates

The execution of a program $\text{syncUpd}(u); P$ remains blocked until the update u is performed over the global store. This is achieved by continuously pulling (i.e., a busy-waiting) until the updates are confirmed by global store. In order to provide the formal semantics of the language, we consider the following runtime syntax for programs.

$$\text{(RUN-TIME-PROGRAM)} \quad P ::= \dots \mid \text{wait}; P \mid e \triangleright (P); P$$

The operational semantics for the new primitives is given by the rules in Fig. 2. Rule (SYNC-UPD) rewrites each synchronous update as the sequence consisting of an asynchronous update followed by `pull` and `wait`. Processes `wait` continuously checks whether local changes have been confirmed by the global store. As described by rule (WAIT), it is implemented as a busy-waiting loop that first checks the local buffers by executing `confirmed` and then performs the conditional jump $x \triangleright (\text{pull}; \text{wait}); P$. If the condition x is true, then it follows as P otherwise it continues as `pull; wait`, as described by rules (GUARD-TRUE) and (GUARD-FALSE).

Single order is characterised, essentially, by imposing arbitration and visibility to coincide [2]. Since our definition for `AR` and `VIS` makes them disjoint, we use an alternative characterisation of single order guarantee, which disregards the arbitration order of updates that are not observed. Hence, we use the following characterisation for single order:

$$\text{AR}; \text{VIS} \subseteq \text{VIS} \quad \text{and} \quad \text{AR}^{-1}; \neg \text{VIS} \subseteq \neg \text{VIS}$$

The following result shows that any well-formed GSP system, whose programs are free from asynchronous updates enjoy the single order guarantee.

Theorem 6.1 (Single Order). *Let N be a well-formed system s.t. $\text{update}(u)$ does not appear in N . If $\langle N, \emptyset, \emptyset, \emptyset, \emptyset \rangle \rightarrow_i^* \langle N', \text{OP}, \text{SS}, \text{SO}, \text{VIS}, \text{AR} \rangle$ then $\text{AR}; \text{VIS} \subseteq \text{VIS}$ and $\text{AR}^{-1}; \neg \text{VIS} \subseteq \neg \text{VIS}$.*

7 Conclusions

We have proposed a formal model for the Global Sequence Protocol and its proposed implementation. We use our formal model to provide a simplified proof (that relies on standard simulation) that the proposed implementation is correct. We remark that our proof does not require to exhibit an auxiliary state for the simulation and that several invariants are trivially ensured by the definition of the model (e.g., the fact that clients have a consistent view of the global sequence) and the well-formed conditions imposed over systems. We have formally studied the consistency guarantees ensured by the model by relying on the operational semantics of the calculus to incrementally compute (a relaxed version of) abstract histories. We have also shown how GSP can be used to formally study programming patterns, like synchronous update operations, that provide stronger consistency guarantees at the expenses of efficiency and availability. We plan to use the GSP calculus as a formal basis for developing programming techniques to enable the fine-tuning of consistency levels in applications.

Acknowledgments We thank the anonymous reviewers of Coordination 2016 for their careful reading of our paper and detailed comments.

References

1. P. Bailis and A. Ghodsi. Eventual consistency today: limitations, extensions, and beyond. *Commun. ACM*, 56(5):55–63, 2013.
2. S. Burckhardt. Principles of eventual consistency. *Foundations and Trends in Programming Languages*, 1(1-2):1–150, 2014.
3. S. Burckhardt, M. Fähndrich, D. Leijen, and B. P. Wood. Cloud types for eventual consistency. In *ECOOP 2012*, pages 283–307. Springer, 2012.
4. S. Burckhardt, D. Leijen, M. Fähndrich, and M. Sagiv. Eventually consistent transactions. In *ESOP 2012*, pages 67–86. Springer, 2012.
5. S. Burckhardt, D. Leijen, J. Protzenko, and M. Fähndrich. Global sequence protocol: A robust abstraction for replicated shared state. In *ECOOP 2015*, pages 568–590, 2015.
6. G. DeCandia, D. Hastorun, M. Jampani, G. Kakulapati, A. Lakshman, A. Pilchin, S. Sivasubramanian, P. Voshall, and W. Vogels. Dynamo: amazon’s highly available key-value store. In *SOSP ’07*, pages 205–220. ACM, 2007.
7. S. Gilbert and N. Lynch. Brewer’s conjecture and the feasibility of consistent, available, partition-tolerant web services. *SIGACT News*, 33(2):51–59, June 2002.
8. A. Gotsman, H. Yang, C. Ferreira, M. Najafzadeh, and M. Shapiro. ‘Cause i’m strong enough: reasoning about consistency choices in distributed systems. In *POPL 2016*, pages 371–384, 2016.
9. A. Lakshman and P. Malik. Cassandra: a decentralized structured storage system. *ACM SIGOPS Operating Systems Review*, 44(2):35–40, 2010.
10. M. Shapiro, N. M. Pregoça, C. Baquero, and M. Zawirski. Conflict-free replicated data types. In *SSS 2011*, pages 386–400, 2011.
11. K. Sivaramakrishnan, G. Kaki, and S. Jagannathan. Declarative programming over eventually consistent data stores. In *PLDI 2015*, pages 413–424. ACM, 2015.