



## Composing ordered sequential consistency



Kfir Lev-Ari<sup>a,\*</sup>, Edward Bortnikov<sup>b</sup>, Idit Keidar<sup>a,b</sup>, Alexander Shraer

<sup>a</sup> Viterbi Department of Electrical Engineering, Technion, Haifa, Israel

<sup>b</sup> Yahoo Research, Haifa, Israel

### ARTICLE INFO

#### Article history:

Received 23 November 2016

Accepted 20 March 2017

Available online 22 March 2017

Communicated by Gregory Chockler

#### Keywords:

Composability

Consistency

Distributed systems

### ABSTRACT

We define *ordered sequential consistency* (OSC), a generic criterion for concurrent objects. We show that OSC encompasses a range of criteria, from sequential consistency to linearizability, and captures the typical behavior of real-world coordination services, such as ZooKeeper. A straightforward composition of OSC objects is not necessarily OSC, e.g., a composition of sequentially consistent objects is not sequentially consistent. We define a global property we call *leading ordered operations*, and prove that it enables correct OSC composition.

© 2017 Elsevier B.V. All rights reserved.

## 1. Introduction

In this work we define a generic correctness criterion named *Ordered Sequential Consistency* (OSC), which captures a range of criteria, from sequential consistency [1] to linearizability [2].

We use OSC to capture the semantics of coordination services such as ZooKeeper [3]. These coordination services provide so-called “strong consistency” for updates and some weaker semantics for reads. They are replicated for high-availability, and each client submits requests to one of the replicas. Reads are not atomic so that they can be served fast, i.e., locally by any of the replicas, whereas update requests are serialized via a quorum-based protocol based on Paxos [4]. Since reads are served locally, they can be somewhat stale but nevertheless represent a valid system state.

In the literature, these services’ guarantees are described as atomic writes and FIFO ordered operations for each client [3]. This definition is not tight in two ways: (1) linearizability of updates has no meaning when no operation reads the written values; and (2) this definition

allows read operations to read from a future write, which obviously does not occur in any real-world service. A special case of OSC, which we call  $OSC(U)$ , captures the actual guarantees of existing coordination services.

Although supporting  $OSC(U)$  semantics instead of atomicity of all operations enables fast local reads, this makes services *non-composable*: correct  $OSC(U)$  coordination services may fail to provide the same level of consistency when combined [5]. Intuitively, the problem arises because  $OSC(U)$ , similarly to sequential consistency [1], allows subset of operations to occur “in the past”, which can introduce cyclic dependencies.

In a companion systems paper [5] we present ZooNet, a system for modular composition of coordination services, which addresses this challenge: Consistency is achieved on the client side by judiciously adding synchronization requests called *leading ordered operations*. The key idea is to place a “barrier” that limits how far in the past reads can be served from. ZooNet does so by adding a “leading” update request prior to a read request whenever the read is addressed to a different service than the previous one accessed by the same client. We provide here the theoretical underpinnings for the algorithm implemented in ZooNet.

Proving the correctness of ZooNet is made possible by the OSC definition that we present in this paper. Interestingly, Vitenberg and Friedman [6] showed that sequential

\* Corresponding author.

E-mail address: kflra@campus.technion.ac.il (K. Lev-Ari).

consistency, when combined with any local (i.e., composable) property continues to be non-composable. Our approach circumvents this impossibility result since having leading ordered operations is not a local property.

## 2. Model and notation

We use a standard shared memory execution model [2], where a set  $\phi$  of sequential processes access shared objects from some set  $X$ . An object has a name label, a value, and a set of operations used for manipulating and reading its value. An operation's execution is delimited by two events, *invoke* and *response*.

A history  $\sigma$  is a sequence of operation *invoke* and *response* events. An *invoke* event of operation  $op$  is denoted  $i_{op}$ , and the matching *response* event is denoted  $r_{op}$ . For two events  $e_1, e_2 \in \sigma$ , we denote  $e_1 <_{\sigma} e_2$  if  $e_1$  precedes  $e_2$  in  $\sigma$ , and  $e_1 \leq_{\sigma} e_2$  if  $e_1 = e_2$  or  $e_1 <_{\sigma} e_2$ . For two operations  $op$  and  $op'$  in  $\sigma$ ,  $op$  precedes  $op'$ , denoted  $op <_{\sigma} op'$ , if  $r_{op} <_{\sigma} i_{op'}$ , and  $op \leq_{\sigma} op'$  if  $op = op'$  or  $op <_{\sigma} op'$ . Two operations are *concurrent* if neither precedes the other.

For a history  $\sigma$ ,  $complete(\sigma)$  is the sequence obtained by removing all operations with no response events from  $\sigma$ . A history is *sequential* if it begins with an *invoke* event and consists of an alternating sequence of *invoke* and *response* events, s.t. each *invoke* is followed by the matching *response*.

For  $p \in \phi$ , the *process subhistory*  $\sigma|p$  of a history  $\sigma$  is the subsequence of  $\sigma$  consisting of events of process  $p$ . The *object subhistory*  $\sigma_x$  for an object  $x \in X$  is similarly defined. A history  $\sigma$  is *well-formed* if for each process  $p \in \phi$ ,  $\sigma|p$  is sequential. For the rest of our discussion, we assume that all histories are well-formed. The order of operations in  $\sigma|p$  is called the *process order* of  $p$ .

For the sake of our analysis, we assume that each subhistory  $\sigma_x$  starts with a dummy initialization of  $x$  that updates it to a dedicated initial value  $v_0$ , denoted  $di_x(v_0)$ , and that there are no concurrent operations with  $di_x(v_0)$  in  $\sigma_x$ .

We refer to an operation that changes the object's value as an *update operation*. The *sequential specification* of an object  $x$  is a set of allowed sequential histories in which all events are associated with  $x$ . For example, the sequential specification of a read-write object is the set of sequential histories in which each read operation returns the value written by the last update operation that precedes it.

## 3. Ordered sequential consistency

**Definition 1** (*OSC(A)*). A history  $\sigma$  is *OSC w.r.t. a subset A of the objects' operations* if there exists a history  $\sigma'$  that can be created by adding zero or more response events to  $\sigma$ , and there is a sequential permutation  $\pi$  of  $complete(\sigma')$ , satisfying the following:

OSC<sub>1</sub> (sequential specification):  $\forall x \in X, \pi_x$  belongs to the sequential specification of  $x$ .

OSC<sub>2</sub> (process order): For two operations  $o$  and  $o'$ , if  $\exists p \in \phi : o <_{\sigma|p} o'$  then  $o <_{\pi} o'$ .

OSC<sub>3</sub> (*A-real-time order*):  $\forall x \in X$ , for an operation  $o \in A$  and an operation  $o'$  (not necessarily in  $A$ ) s.t.  $o, o' \in \sigma_x$ , if  $o' <_{\sigma} o$  then  $o' <_{\pi} o$ .

Such  $\pi$  is called a *serialization* of  $\sigma$ . An object is *OSC(A)* if all of its histories are *OSC(A)*.

We assume that  $\forall x \in X, di_x(v_0) \in A$ . Linearizability and sequential consistency are both special cases of *OSC(A)*: (1) we get linearizability using  $A$  that consist of all of the objects' operations; and (2) we get sequential consistency with  $A$  that consists only of dummy initialization operations, which means that there is no operation that precedes an  $A$ -operation, i.e., OSC<sub>3</sub> is null, and we left with the sequential specification and process order of an object.

If  $A$  consists of the objects' update operations, denoted  $U$ , then *OSC(U)* captures the semantics of coordination services: (1) updates are globally ordered (by OSC<sub>3</sub>); and (2) all operations see some prefix of that order (by OSC<sub>3</sub>), while respecting each client process order (by OSC<sub>2</sub>).

## 4. OSC(A) composability via leading A-operations

In this section we show that a history  $\sigma$  of *OSC(A)* objects satisfies *OSC(A)*, if  $\sigma$  has leading ordered  $A$ -operations. Generally, we prove the composition by ordering every  $A$ -operation  $o_A$  on object  $x$ , according to the first event  $e \in \sigma$  s.t.  $e \leq_{\sigma} r_{o_A}$  and  $i_{o_A} <_{\pi_x} e$ . Then, we extend that order to a total order on all operations, by placing every non- $A$ -operation after the  $A$ -operation that precedes it in their object's serialization. Finally, we show that if  $\sigma$  has leading ordered  $A$ -operations, then the total order satisfies *OSC(A)*. Intuitively, we can think of the leading  $A$ -operations as a barrier for the non- $A$ -operations, that maintains the total order between objects.

Given a history  $\sigma$  of *OSC(A)* objects, and a set of serializations  $\Pi = \{\pi_x\}_{x \in X}$  of  $\{\sigma_x\}_{x \in X}$ , we define a strict total order on all operations in  $\Pi$ . We refer to an operation  $o \in A$  as an  $A$ -operation, and define the future set of an  $A$ -operation as follows:

**Definition 2** (*A-operation future set*). Given a history  $\sigma$  of *OSC(A)* objects, an object  $x \in \sigma$ , a serialization  $\pi_x$  of  $\sigma_x$ , and an  $A$ -operation  $o_A \in \sigma_x$ , the *future set of  $o_A$  in  $\pi_x$*  is  $F_{\sigma}^{\pi_x}(o_A) \triangleq \{o \in \pi_x | o_A \leq_{\pi_x} o\}$ .

We now define an  $A$ -operation's first response event to be the earliest response event of an operation in its future set.

**Definition 3** (*First response event*). Given a history  $\sigma$  of *OSC(A)* objects, an object  $x \in \sigma$ , a serialization  $\pi_x$  of  $\sigma_x$ , and an  $A$ -operation  $o_A \in \pi_x$ , the *first response event of  $o_A$  in  $\pi_x$* , denoted  $fr_{\sigma}^{\pi_x}(o_A)$ , is the earliest response event in  $\sigma$  of an operation in  $F_{\sigma}^{\pi_x}(o_A)$ .

Note that it is possible that  $fr_{\sigma}^{\pi_x}(o_A)$  is  $o_A$ 's response event. We make two observations regarding first responses:

Download English Version:

<https://daneshyari.com/en/article/4950785>

Download Persian Version:

<https://daneshyari.com/article/4950785>

[Daneshyari.com](https://daneshyari.com)