

# Fast linearizability testing

Alex Horn and Daniel Kroening

University of Oxford

February 5, 2015

## Abstract

Linearizability is the de facto correctness criteria for concurrent data structures and consensus-based distributed systems. In this talk, I present some work-in-progress on a new tool for testing linearizability. We experimentally show that our tool can handle problems where existing implementations timeout or run out of memory. Our experiments include shared memory programs (Intel's TBB and Siemens' EMBB library) and a large-scale distributed system (etcd) using Aphyr's testing framework (Jepsen).

# Outline

## What's the problem?

- Example

- Problem statement

## Ideas

- Partitioning scheme

- Optimizations

## Experiments

- Etcd distributed system

- Intel's TBB library

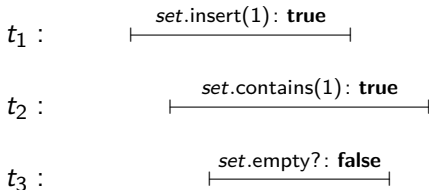
- Siemens' EMBB library

## Concluding remarks

## Example: history

Consider a concurrent data type for unordered sets with the usual 'insert', 'contains' and 'empty?' operations.

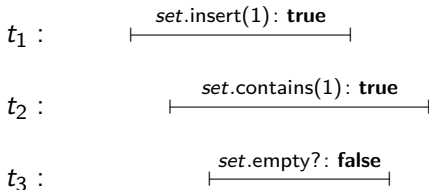
Suppose three threads execute these operations concurrently, and we record their operations in a so-called **history** (let it be  $H_1$ ):



What does it mean for this execution to be correct (or incorrect)?  
This is answered by the concept of **linearizability** (see next slides).

## Example: linearizability

Let  $H_1$  be the history from the previous slide:



Then  $H_1$  is linearizable because we can reorder the operations to, say,

*i.* `set.insert(1): true`; *ii.* `set.empty?: false`; *iii.* `set.contains(1): true`

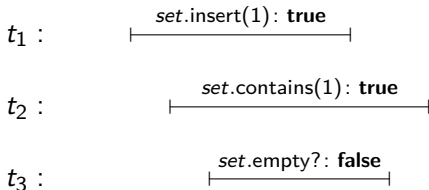
which satisfies the specification of a sequential set.

The problem (to be formalized next) is to automatically find such a reordering given any  $H$  and sequential specification.

## Definition (History)

Let  $E \triangleq \{\text{call}, \text{ret}\} \times \mathbb{N}$  be the set of **events**. For all  $n$  in  $\mathbb{N}$ , define  $\text{call}_n \in E$  to be a **call** and  $\text{ret}_n \in E$  to be a **return**. A **history**  $H$  is a sequence of such events, totally ordered by  $\preceq_H$ . For all events  $e$  in  $E$ ,  $\text{obj}(e)$  and  $m(e)$  denotes the **object** and **method** of  $e$ , respectively. The **length** of  $H$ , written  $|H|$ , is the number of events in  $H$ .

## Example



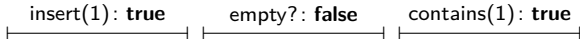
$H_1 \triangleq \langle \text{call}_1, \text{call}_2, \text{call}_3, \text{ret}_1, \text{ret}_3, \text{ret}_2 \rangle$  where  $\text{obj}(\text{call}_1) = \text{'set'}$  and  $m(\text{call}_1) = \text{'insert}(1): \text{true}'}$ ,  $|H_1| = 6$  and  $\text{call}_3 \preceq_{H_1} \text{ret}_1$  etc.

## Definition (Sequential history)

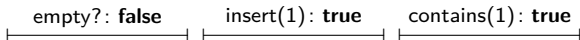
Let  $e, e' \in E$  be events in history  $H$ . If  $e$  is a call and  $e'$  is a return in  $H$ , both are **matching** whenever  $e \preceq_H e'$  and their objects and operations are equal, i.e.  $obj(e) = obj(e')$  and  $m(e) = m(e')$ . A history is **complete** if every call has a matching return. A complete history  $H$  is **sequential** if it alternates between matching calls and returns (necessarily starting with a call).

## Example

$H_2$  is a sequential history:



And so is  $H_3$  (albeit not maybe what we expect):



## Definition (Happens-before)

Given a history  $H$ , the **happens-before** relation is defined to be a partial order  $<_H$  over calls  $e$  and  $e'$  such that  $e <_H e'$  whenever  $\text{ret}(e) \preceq_H e'$ , i.e.  $e$ 's matching return occurs before  $e'$  in  $H$ .

## Example

Let  $H_5$  be the following history:

$t_1$  :     $\text{call}_1 \xrightarrow{\text{set.insert}(1): \text{true}} \text{ret}_1$

$t_2$  :     $\text{call}_3 \xrightarrow{\text{set.contains}(1): \text{true}} \text{ret}_3$

$t_3$  :     $\text{call}_2 \xrightarrow{\text{set.empty?}: \text{false}} \text{ret}_2$

- $\text{call}_1 \not<_H \text{call}_2$  and  $\text{call}_2 \not<_H \text{call}_1$ , i.e. both **happen concurrently**;
- $\text{call}_1 <_H \text{call}_3$  and  $\text{call}_2 <_H \text{call}_3$ , i.e.  $\text{call}_1$  **happens-before**  $\text{call}_3$  etc.



# Problem statement

## Definition (Linearizability)

Denote with  $\phi$  a specification of a sequential data type, i.e.  $\phi$  is a unary predicate on sequential histories. A  $\phi$ -**sequential history** is a sequential history  $H$  that satisfies  $\phi(H)$ . A history  $H$  is **linearizable with respect to  $\phi$**  if it can be extended to a complete history  $H'$  (by appending zero or more returns to  $H$ ) and there is a  $\phi$ -sequential history  $S$  such that

**L1**  $H'$  is equivalent to  $S$ , i.e.  $H'$  and  $S$  are equal as sets;

**L2**  $<_H \subseteq <_S$ , i.e. if  $e$  precedes  $e'$  in  $H$ , the same is true in  $S$ .

**Decision problem:** Given  $H$  and  $\phi$ , decide whether  $H$  is linearizable with respect to  $\phi$ .

## Example: linearizability (with our formalization)

Recall  $H_5$ :

$t_1$  :     $\text{call}_1 \xrightarrow{\text{set.insert}(1): \text{true}} \text{ret}_1$

$t_2$  :     $\text{call}_3 \xrightarrow{\text{set.contains}(1): \text{true}} \text{ret}_3$

$t_3$  :     $\text{call}_2 \xrightarrow{\text{set.empty?}: \text{false}} \text{ret}_2$

**Question:** Is  $H_5$  linearizable with respect to  $\phi_{\text{set}}$ ?

## Example: linearizability (with our formalization)

Recall  $H_5$ :

$t_1$ :     $\text{call}_1 \xrightarrow{\text{set.insert}(1): \text{true}} \text{ret}_1$

$t_2$ :     $\text{call}_3 \xrightarrow{\text{set.contains}(1): \text{true}} \text{ret}_3$

$t_3$ :     $\text{call}_2 \xrightarrow{\text{set.empty?}: \text{false}} \text{ret}_2$

**Question:** Is  $H_5$  linearizable with respect to  $\phi_{\text{set}}$ ? Yes!

$H_2$  is a witness for a  $\phi_{\text{set}}$ -sequential history that respects  $<_{H_5}$ :

$\xrightarrow{\text{insert}(1): \text{true}} \xrightarrow{\text{empty?}: \text{false}} \xrightarrow{\text{contains}(1): \text{true}}$

In general, to check linearizability, we need to consider all permutations of a history  $H$ . That's  $O(|H|!)$ .

## Theorem (Complexity)

*Given a history  $H$  and specification  $\phi$ , the linearizability problem — i.e. whether  $H$  is linearizable with respect to  $\phi$  — is NP-complete.*

### Proof.

The problem is clearly in NP whenever the specification of the sequential data type is deterministic and can be run in polynomial time. For NP-hardness, by reduction from SAT. Let  $\psi$  be a propositional logic formula. Let  $M_\psi$  be a model for a partial Boolean assignment of the variables in  $\psi$ . Let  $H$  be a history where every operation sets a variable to either true or false, and all operations happen concurrently, i.e. are unordered by  $<_H$ . The specification is satisfied whenever all variables in  $M_\psi$  have been assigned a value without causing a contradiction (any later assignments are ignored). Then  $H$  is linearizable with respect to  $M_\psi$  if and only if  $\psi$  is SAT. Hence linearizability checking is NP-complete. □

## About our decision procedure

### Definition (Interval order)

A partially ordered set  $(P, \leq)$  is an **interval order** whenever, for all  $x, y, u, v$  in  $P$ , if  $x \leq y$  and  $u \leq v$ , then  $x \leq v$  or  $u \leq y$ .

Note that  $(P, \leq)$  is an interval order if and only if no restriction of  $(P, \leq)$  is isomorphic to the following Hasse diagram:<sup>1</sup>



---

<sup>1</sup>Rabinovitch, I., *The dimension of semiorders*. J. Comb. Theory (Series A) 25 1978 50–61

## About our decision procedure

### Theorem

*For every complete history  $H$ ,  $<_H$  is an interval order.*

### Proof.

Assume  $x <_H y$  and  $u <_H v$ . Let  $\text{ret}(x)$  and  $\text{ret}(u)$  be the matching return of  $x$  and  $u$ , respectively. Since  $H$  is complete, both returns are in  $H$ . Recall that  $\preceq_H$  denotes the total ordering of events in  $H$ . By definition of  $<_H$ ,  $\text{ret}(x) \preceq_H y$  and  $\text{ret}(u) \preceq_H v$ . Since  $\preceq_H$  is total, either  $\text{ret}(x) \preceq_H \text{ret}(u)$  or  $\text{ret}(u) \preceq_H \text{ret}(x)$ . The former implies  $x <_H v$ , whereas the latter implies  $u <_H y$ , proving the claim.  $\square$

Put differently, this talk is about a **decision procedure for a certain class of partial orders**.

## Idea: Partitioning scheme

### Theorem

*A history  $H$  is linearizable with respect to a specification if and only if, for each object  $obj$ , the restriction of  $H$  to  $obj$  is linearizable.<sup>2</sup>*

We use this theorem to partition a history of operations on an associative data type (e.g. a set or hash map) into sub-histories that only contain operations on the same datum. We call this the **compositional technique**. Note that this technique cannot be applied when there are zero-arg operations such as 'empty?'.  

---

### Example

`set.insert(1)` and `set.insert(2)` are in different sub-histories even though both operations are on the same set. Put differently, we interpret the datum to be *obj*.

<sup>2</sup>This is well-known as “compositionality”, see Theorem 3.6.1 in *The Art of Multiprocessor Programming* (Revised Ed.) by Herlihy and Shavit.

## Idea: Optimizations

As in SAT solvers, optimizations matter. We've experimentally evaluated many different ways of implementing the decision procedure, including non-chronological backtracking, parallelization, and cache eviction strategies. We also looked at hashing: we exploit the fact that XOR forms an abelian group and by this we get a constant-time,  $O(1)$ , hash function on bitsets.

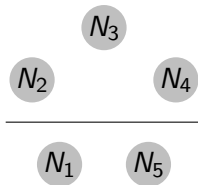
Abelian group axioms:

- Closure: for all  $x, y \in \mathbb{V}$ ,  $x \oplus y \in \mathbb{V}$ ;
- Associativity: for all  $x, y, z \in \mathbb{V}$ ,  $(x \oplus y) \oplus z = x \oplus (y \oplus z)$ ;
- Identity: for all  $x \in \mathbb{V}$ ,  $x \oplus 0 = x = 0 \oplus x$ ;
- Inverse: for all  $x \in \mathbb{V}$ ,  $x \oplus \tilde{x} = 0$  (here,  $\tilde{x} \triangleq x$ );
- Commutativity: for all  $x, y \in \mathbb{V}$ ,  $x \oplus y = y \oplus x$ .



## Experiments with distributed systems

For our experiments with distributed systems, we use *Jepsen*, a testing framework that simulates network failures.<sup>3</sup> For example, Jepsen sporadically introduces network partitions in a cluster (on purpose):



Using Jepsen, we have collected histories from *etcd*, a distributed key/value data store.<sup>4</sup> Due to network failures, the *etcd* histories are rather different from those generated by shared memory programs.

---

<sup>3</sup><https://github.com/aphyr/jepsen>

<sup>4</sup><https://github.com/coreos/etcd>

## Etcd experiments, comparison with Knossos

Jepsen ships with a linearizability tester called Knossos.<sup>5</sup> Out of 100 etcd histories collected with Jepsen, 80% are non-linearizable (this is harder to check than linearizability, i.e. SAT versus UNSAT).

Crucially, Knossos times out on benchmarks 7 and 99, and runs out of memory on benchmarks 40, 57, 85 and 97. All these benchmarks (except 7) are non-linearizable. In contrast, our linearizability tester completes all 100 benchmarks in around 7 s, averaging **0.07 s** per benchmark. Memory usage of our tool is negligible.

---

<sup>5</sup><https://github.com/aphyr/knossos>

## TBB experiments

We tested `concurrent_unordered_set` in Intel's *TBB* library.<sup>6</sup> In our experiments, we start up four processes that randomly insert and remove items. Each process produces 70 K operations. In other words, the length of histories is  $4 \times 2 \times 70K = 560 K$ . Those histories are significantly longer than in previous experiments where the limit is  $2^{13} \approx 8 K$  operations per process.<sup>7</sup> For our tool, we get the following results when histories are linearizable:

- LRU=on, O(1) hash: **54 s [646 MiB]**
- LRU=off, O(1) hash: 131 s [9,763 MiB  $\approx$  10 GiB]
- Compositional: **7 s [9,763 MiB  $\approx$  10 GiB]**

The first two times remain stable even if allow 'empty?' checks, but then we cannot use the compositional technique (as noted before).

---

<sup>6</sup><https://www.threadingbuildingblocks.org/>

<sup>7</sup><http://www.cs.ox.ac.uk/people/gavin.lowe/LinearizabilityTesting/>

## EMBB experiments

We tested LockFreeStack in Siemens' EMBB library.<sup>8</sup> Again, we start up four processes that randomly push and pop items. Each process produces 70 K operations, resulting in histories of length 560 K. We get the following results when histories are linearizable:

- LRU=on, O(1) hash: 73 s [**704 MiB**]
- LRU=off, O(1) hash: **59 s** [13, 135 MiB  $\approx$  13GiB]

Things to improve: Currently, we artificially induce non-linearizable histories by changing the initial state. Ideally, we would find real bugs in TBB or EMBB. It would be also interesting to extract histories from Gavin's experiments on hash sets implemented in Scala.<sup>9</sup>

---

<sup>8</sup><https://github.com/siemens/embb>

<sup>9</sup><http://www.cs.ox.ac.uk/people/gavin.lowe/LinearizabilityTesting/>

## Concluding remarks

We have presented a linearizability tester with promising results. We beat Knossos and our partitioning scheme shows significant algorithmic improvements over Gavin's linearizability tester. In addition, we have shown that LRU cache eviction can reduce memory consumption by one order of magnitude. In fact, in the TBB experiments, this even decreases runtime.

As future work, it would be nice to extend the idea behind our partitioning scheme to a wider class of operations. I am planning to support other consistency models.

You may download our tool now!

<https://github.com/ahorn/linearizability-tester>

(Our source code repository contains not only our experiments but also many unit tests to simplify development.)