Consistency Criteria for Distributed Data Structures

Constantin Enea IRIF, University Paris Diderot

joint work with Ahmed Bouajjani, Rachid Guerraoui, Jad Hamza

Ĭ



European Research Council Established by the European Commission



Geo-Distributed Software Systems



- large numbers of users
- distributed across wide geographical areas
- generate and access huge amounts of data

Examples: - online markets (Amazon)

- government services (tax payment)
- enterprise customer services
- social networks (Facebook)
- managing connected devices and sensors (Internet of Things)

Distributed Data Storage (Distributed Data Structures)



- to support failures, data is **replicated**
- for availability, replicas may store different versions of data: weak consistency
- interface restricted to a **fixed set of methods**
 - key-value stores provide put(key,value) and get(key)
- weak consistency and restricted interface: different approach than relational databases

Current Status

Data storage layer

- informal and vague specifications
- lack of rigorous methodologies for validation
 - testing the behavior of the system under stress
- dormant bugs with potential severe consequences (millions of users)
 - used in critical systems: developing drugs, managing medical equipment, etc.

Application layer

- lack of programming abstractions of data storage while developing applications
 - no compositional reasoning
- complex and fragile software systems

Optimistic replication: replicas are allowed to diverge

• operations are applied immediately at the submission site



Optimistic replication: replicas are allowed to diverge

• operations are applied immediately at the submission site



Optimistic replication: replicas are allowed to diverge

- operations are applied immediately at the submission site
- in the background, sites exchange and apply remote operations



Optimistic replication: replicas are allowed to diverge

- operations are applied immediately at the submission site
- in the background, sites exchange and apply remote operations



Optimistic replication: replicas are allowed to diverge

- operations are applied immediately at the submission site
- in the background, sites exchange and apply remote operations







Solving conflicts between concurrent operations



Solving conflicts between concurrent operations



Solving conflicts between concurrent operations



Solving conflicts between concurrent operations



Solving conflicts between concurrent operations



Solving conflicts between concurrent operations

- speculate and roll-back, e.g., Google App Engine Datastore
- convergent conflict resolution, e.g., CRDTs [Shapiro et al.'11]



Formal verification of DDSs

- Correct operations ? Allowed level of consistency between replicas ?
 - by CAP theorem [Gilbert et al.'02], achieving strong consistency (linearizability) is impossible
 - various correctness criteria: eventual consistency, causal consistency, etc.
- Developing algorithmic methods for the verification of these criteria

- Formal definitions of consistency criteria
 - safety: each operation is executed in the context of a local view, which must satisfy some specification
 - sites can have different local views
 - liveness: the local views converge toward a global view
- Automatic verification of consistency criteria
 - decidability/complexity results
 - general reductions to classical verification problems (reachability, model checking)

Modeling DDS behaviors with traces

- operations are instances of a set of methods (add, rem, lookup)
- traces record the submitted operations and their return values
 - trace = a partially-ordered set of operations
 - operations submitted to the same site are ordered



• Specifying local views using relations between operations



- Vis(o',o): o' is visible when o is executed
- Order_o(o₁,o₂): when o is executed, o1 has been executed before o2

• Specifying local views using relations between operations



- Vis(o',o): o' is visible when o is executed
- Order_o(0₁,0₂): when o is executed, o1 has been executed before o2

• Specifying local views using relations between operations



- Vis(o',o): o' is visible when o is executed
- Order_o(o₁,o₂): when o is executed, o1 has been executed before o2
- Safety: ∃ Vis ∀o ∃ Order₀ such that ...

- Vis(o',o): o' is visible when o is executed
- Order_o(o₁,o₂): when o is executed, o1 has been executed before o2
- Safety: \exists Vis $\forall \circ \exists$ Order_o such that
 - eventual consistency: Vis U Site_order is acyclic

Ordero satisfies the specification

- Vis(o',o): o' is visible when o is executed
- Order_o(o₁,o₂): when o is executed, o1 has been executed before o2
- Safety: \exists Vis $\forall \circ \exists$ Order_o such that
 - eventual consistency: Vis U Site_order is acyclic

Order_o satisfies the specification

read-your-own-writes: + Site_order ⊆ Vis

if Vis(o₁,o), Vis(o₂,o), and Site_order(o₁,o₂) then Order_o(o₁,o₂)

- Vis(o',o): o' is visible when o is executed
- Order_o(o₁,o₂): when o is executed, o1 has been executed before o2
- Safety: \exists Vis $\forall \circ \exists$ Order_o such that
 - eventual consistency: Vis U Site_order is acyclic

Order_o satisfies the specification

read-your-own-writes: + Site_order ⊆ Vis

if Vis(o₁,o), Vis(o₂,o), and Site_order(o₁,o₂) then Order_o(o₁,o₂)

causal consistency: + Site_order ⊆ Vis and Vis is transitive

if Vis(o₁,o), Vis(o₂,o), and Vis(o₁,o₂) then Order_o(o₁,o₂)

Specifications

- Sequential world: an operation is considered correct depending on the sequence of previously executed operations
 - lookup(0) ▷ true is correct iff it executes after a sequence of operations where the projection on operations with input 0 ends in add(0)
 - e.g., add(2), rem(0), ..., add(0), rem(2)
- **Distributed** world: sequences are replaced by partial orders
 - lookup(0) > true is correct iff it executes after a poset of operations where the projection on operations with input 0 contains a maximal add(0)
 - e.g., add(2), rem(0), ..., add(0), rem(2)

- Arb(01,02): eventually, all sites agree that 01 should be executed before 02
- Safety + Liveness: ∃ Vis ∃ Arb ∀o ∃ Order₀ such that ... and

∀ 01,02: #0 s.t. **Order₀(01,02)** ∧ **Arb(02,01)** is finite

Verification problems

 Given a finite-state implementation Impl and a "regular" specification Spec, is there a "simple" centralized monitor that

outputs error iff Impl is not X-consistent w.r.t. Spec

(X = eventually, causal, read-your-own-writes,...)

- Eventual consistency: Yes (counting + Presburger assertions)
- Read-your-own-writes, Causal consistency: No
- Linearizability: Yes (finite-state)
- Restricting specifications: for causal consistency, there exists such a monitor if the Spec is key-value store

Conclusions

- **Distributed data structures**, an alternative to classical databases
- Several **consistency criteria** which are poorly understood
 - guarantees for the application layer ?
- Hard (undecidable) verification problems
 - finding reasonable restrictions or approximations