

TARDiS: A branch and merge approach to weak consistency

By: Natacha Crooks, Youer Pu, Nancy Estrada, Trinabh Gupta, Lorenzo Alvis, Allen Clement

Presented by: Samodya Abeysiriwardane

TARDiS

Transactional key-value store for weakly consistent systems



Weakly consistent systems

ALPS (Available, low Latency, Partition tolerance, high Scalability)

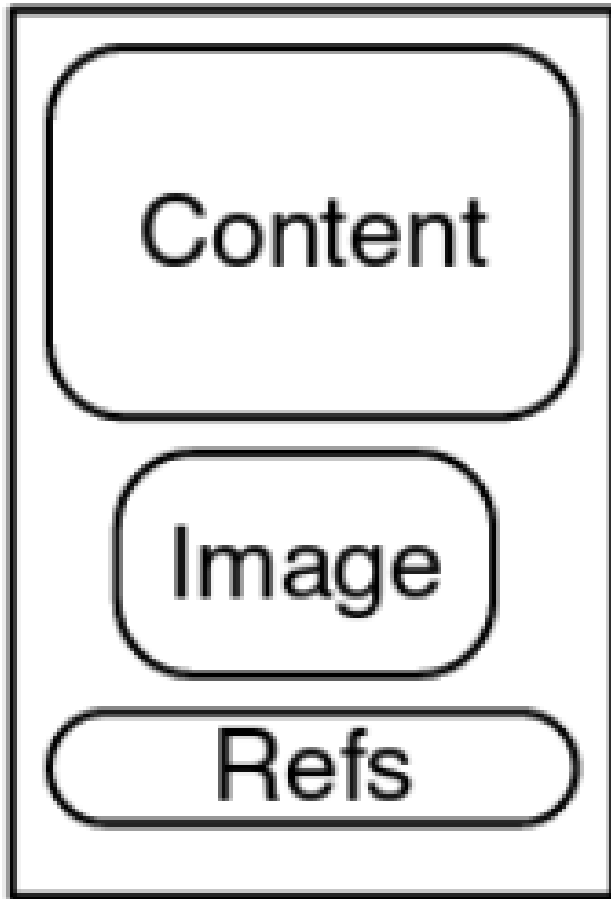
Conflicting operations cause replicas to diverge

Current solutions: Deterministic Writer Wins, per object eventual convergence (object as unit of merging)

Current solutions are not sufficient



Motivation



A wiki page with three objects

Edited at two georeplicated replicas

Motivation

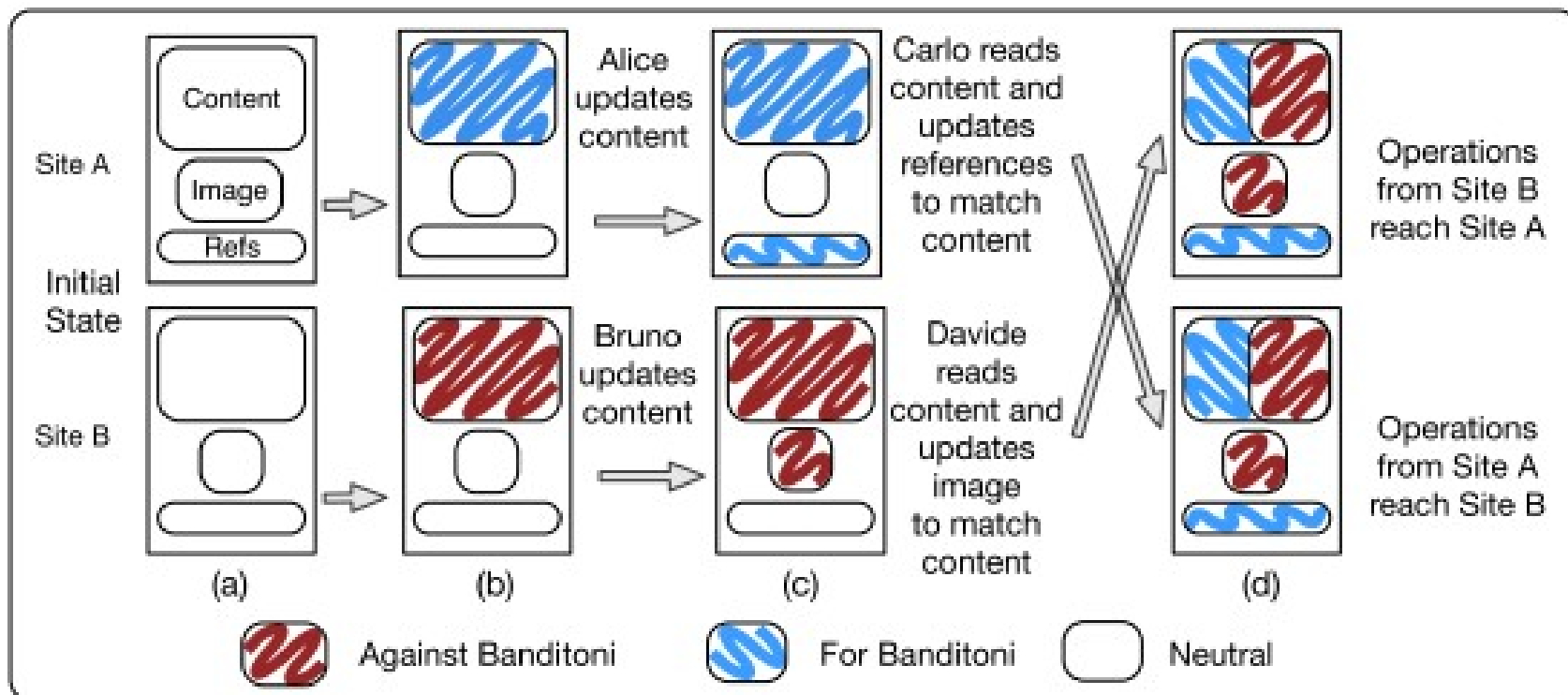


Figure 1: Weakly-consistent Wikipedia

Main goal

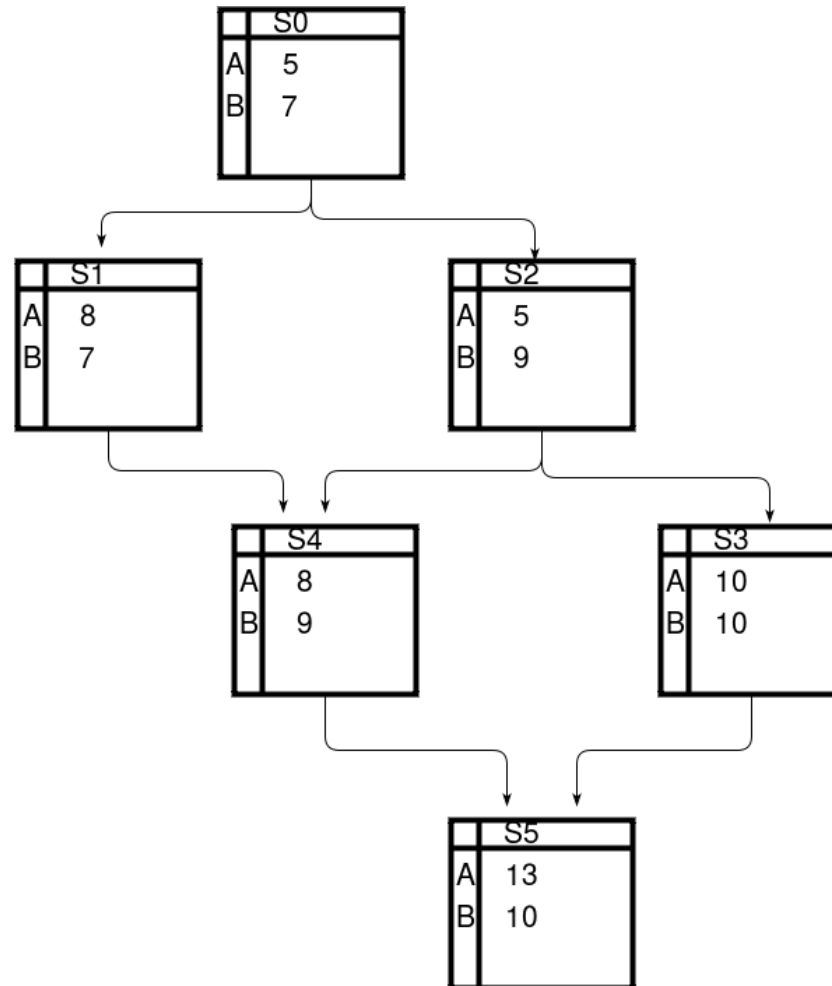
Give applications access to context that is essential for reasoning about concurrent updates

Proposed solution

Expose branches as a unit of merging

- **branch on conflict**
- **branch isolation**
- **application driven merges**

Simple Example with Counters



Key value store
of Counters

Merge

Need to define a merge function for the application

Merging two counters A and B

For counters 2-way merge

```
fn merge (lca, a, b) = lca + (a-lca) + (b-lca)
```

For counters n-way merge

```
fn merge {
```

```
  lca = find_fork_point
```

```
  val = lca
```

```
  for v in conflicting_values:
```

```
    val += (a - lca) + (b - lca)
```

```
}
```

Simple Example with Counter (Code)

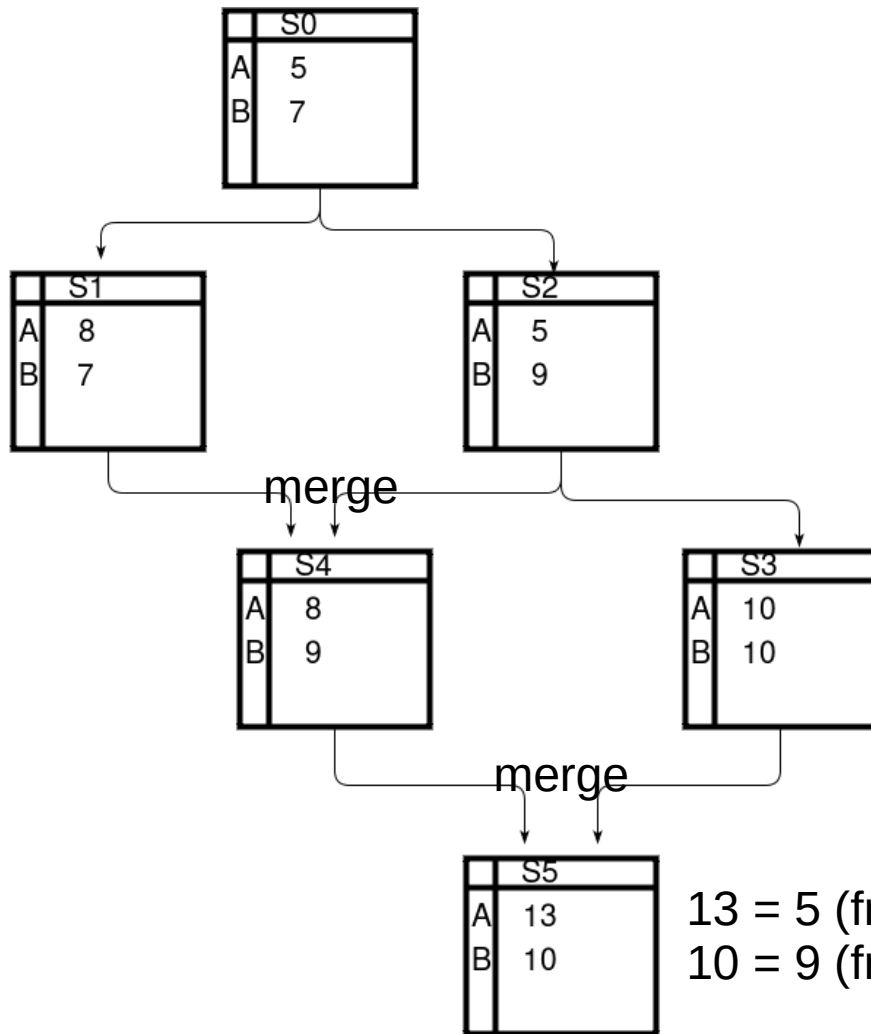
```
1  func increment(counter)
2    Tx t = begin(AncestorConstraint)
3    int value = t.get(counter)
4    t.put(counter, value + 1)
5    t.commit(SerializabilityConstraint)

7  func decrement(counter)
8    Tx t = begin(AncestorConstraint)
9    int value = t.get(counter)
10   t.put(counter, value - 1)
11   t.commit(SerializabilityConstraint)

13 func merge()
14   Tx t = beginMerge(AnyConstraint)
15   forkPoint forkPt =
16     t.findForkPoints(t.parents).first
17   int forkVal = t.getForID(counter, forkPt)
18   list<int> currentVals =
19     t.getForID(counter, t.parents)
20   int result = forkVal
21   foreach c in currentVals
22     result += (c - forkVal)
23   t.put(counter, result)
24   t.commit(SerializabilityConstraint)
```

Figure 3: TARDiS' counter implementation

Simple Example with Counter (Code)



Client1
T1:
inc(A, 3)

Client2
T2:
inc(B,2)

Tm:
merge

T3:
inc(A,5)
inc(B,1)

Tm:
merge

$$13 = 5 \text{ (from S2)} + (8-5) + (10-5)$$
$$10 = 9 \text{ (from S2)} + (9-9) + (10-9)$$

Example

Impose an application invariant of

- if $A > 8$: B should max at 10**

- the merge function can be changed to reflect that**

Highlights the need for cross object merging semantics vs per object merging

Therefore branches as a unit of merging

Another example: Inventory

XYZ_stock: 1

ABC_stock: 3

Alice buys XYZ

XYZ_stock: 0

Bob buys XYZ and ABC

XYZ_stock: 0

ABC_stock: 2

Merge

Bob get XYZ, and exp

Alice gets error

XYZ_stock: 0

exp_stock: 2

Invariant: stock cannot be < 0

Other advantages

No locking required

Branching as a fundamental abstraction for modeling conflicts end to end - replicas as well the local site can be viewed as branches

TARDIS API

S	M	return type	method
✓		transaction	begin(<i>beginConstraint</i>)
	✓	transaction	beginMerge(<i>beginConstraint</i>)
✓	✓	void	put(<i>key</i> , <i>value</i>)
✓		value	get(<i>key</i>)
	✓	value	getForID(<i>key</i> , <i>StateID</i> [])
	✓	key[]	findConflictWrites(<i>StateID</i> [])
	✓	forkPoints[]	findForkPoints(<i>StateID</i> [])
✓	✓	abort commit	commit(<i>endConstraint</i>)

Table 2: TARDiS API - S:single mode, M:merge mode

TARDiS architecture

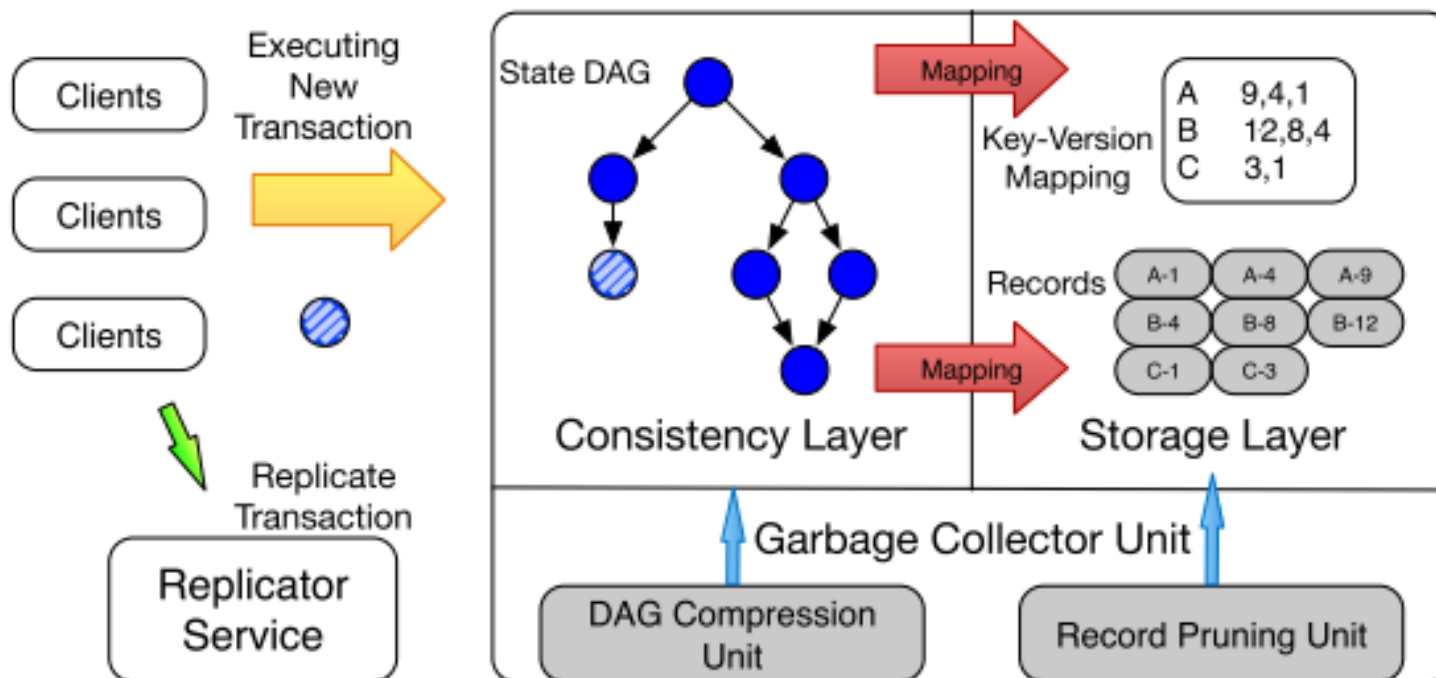


Figure 2: TARDiS architecture

TARDiS architecture

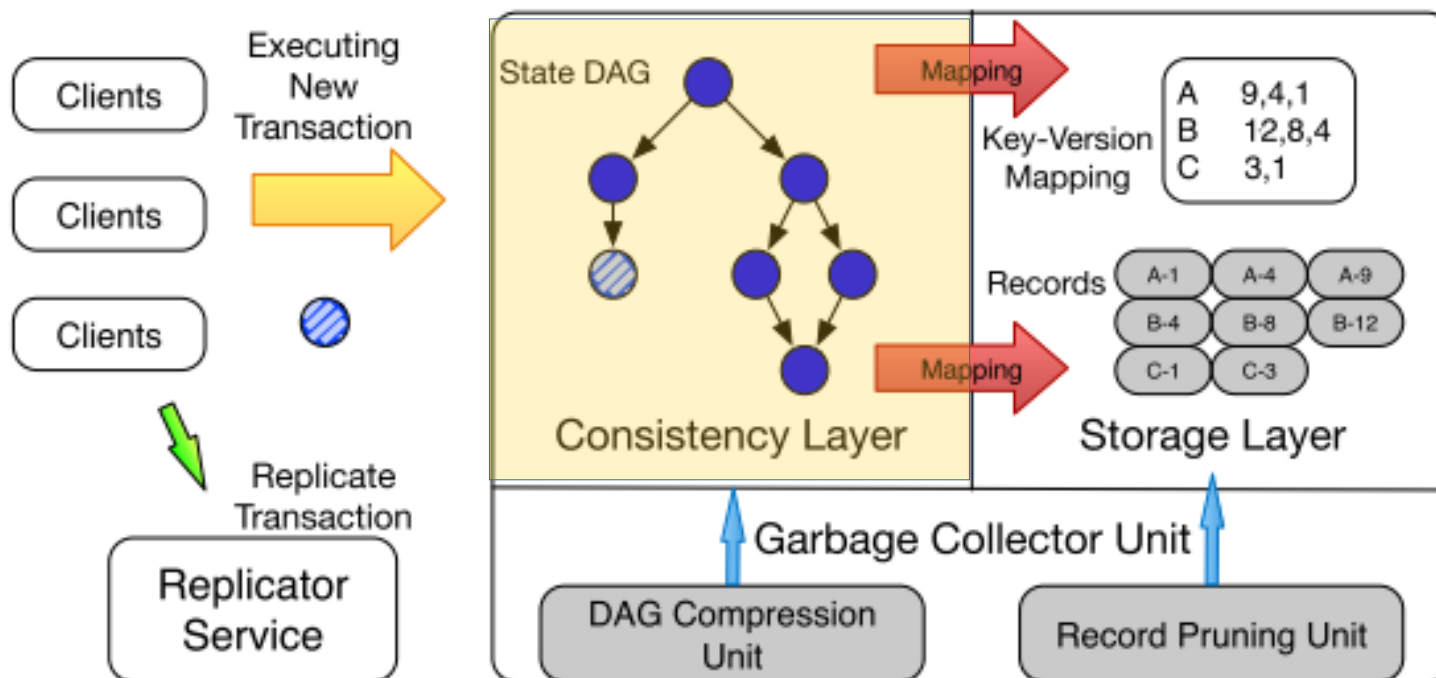


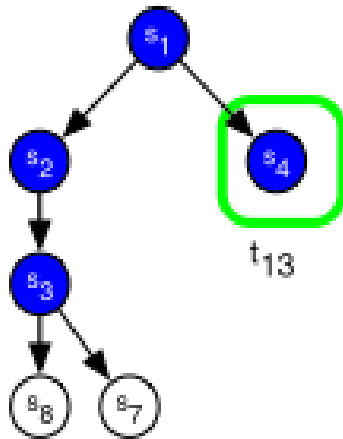
Figure 2: TARDiS architecture

Consistency layer

Constraint	B	E	Description
<i>Any</i>	✓	✓	Always Satisfies
<i>Serializability</i>		✓	Guarantees Serializability
<i>Snapshot Iso</i>		✓	Guarantees Snapshot Isolation
<i>Read Committed</i>		✓	Guarantees Read Committed
<i>No Branching</i>		✓	State has no children
<i>K-Branching</i>		✓	State has fewer than k-1 children
<i>Parent</i>	✓		State where client last committed
<i>Ancestor</i>	✓		Child of client's last committed state
<i>State Identifier</i>	✓		State ID matches the specified ID

Table 1: Begin (B) and end (E) constraints supported by TARDiS

Consistency layer



begin(AncestorConstraint)



Begin & end
constraint unsatisfied



Begin & end
constraint satisfied



Final Committed State



Selected read state



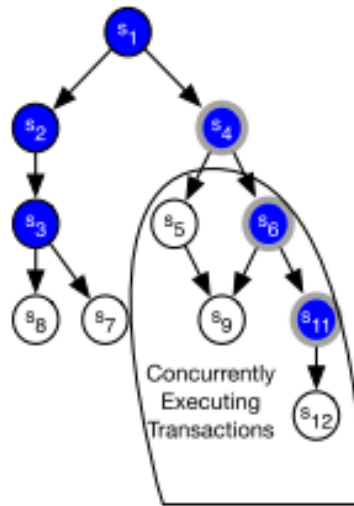
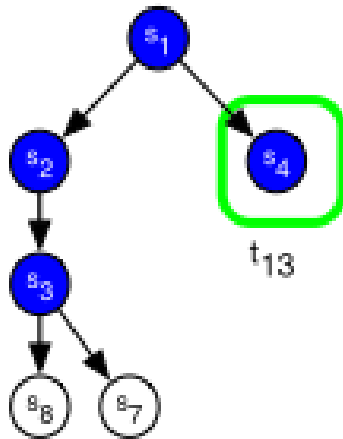
Candidate commit state

Consistency layer

Constraint	B	E	Description
<i>Any</i>	✓	✓	Always Satisfies
<i>Serializability</i>		✓	Guarantees Serializability
<i>Snapshot Iso</i>		✓	Guarantees Snapshot Isolation
<i>Read Committed</i>		✓	Guarantees Read Committed
<i>No Branching</i>		✓	State has no children
<i>K-Branching</i>		✓	State has fewer than k-1 children
<i>Parent</i>	✓		State where client last committed
<i>Ancestor</i>	✓		Child of client's last committed state
<i>State Identifier</i>	✓		State ID matches the specified ID

Table 1: Begin (B) and end (E) constraints supported by TARDiS

Consistency layer



Begin & end
constraint unsatisfied



Selected read state



Begin & end
constraint satisfied



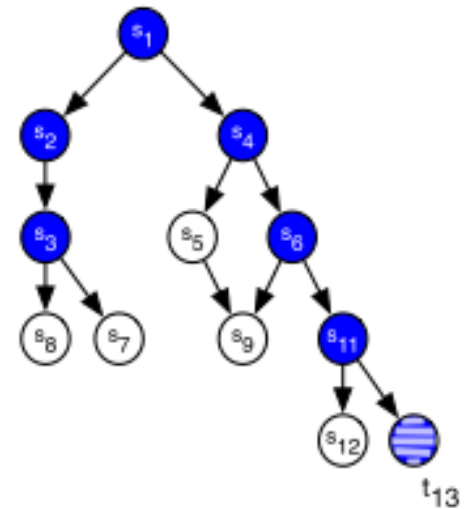
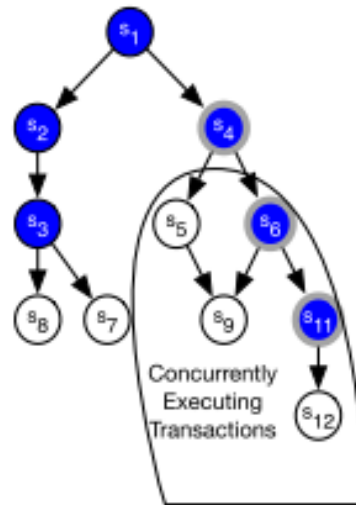
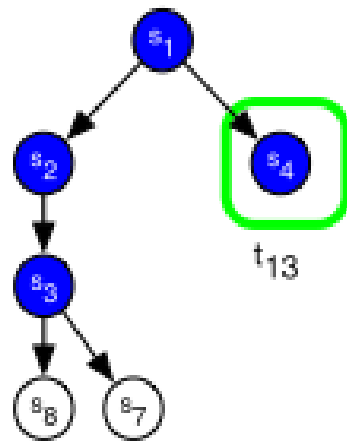
Candidate commit state



Final Committed State

Consistency layer

commit(SerializabilityConstraint)



○ Begin & end
constraint unsatisfied

● Begin & end
constraint satisfied

⊘ Final Committed State

□ Selected read state

⦿ Candidate commit state

Consistency layer

Constraint	B	E	Description
<i>Any</i>	✓	✓	Always Satisfies
<i>Serializability</i>		✓	Guarantees Serializability
<i>Snapshot Iso</i>		✓	Guarantees Snapshot Isolation
<i>Read Committed</i>		✓	Guarantees Read Committed
<i>No Branching</i>		✓	State has no children
<i>K-Branching</i>		✓	State has fewer than k-1 children
<i>Parent</i>	✓		State where client last committed
<i>Ancestor</i>	✓		Child of client's last committed state
<i>State Identifier</i>	✓		State ID matches the specified ID

Table 1: Begin (B) and end (E) constraints supported by TARDiS

TARDiS architecture

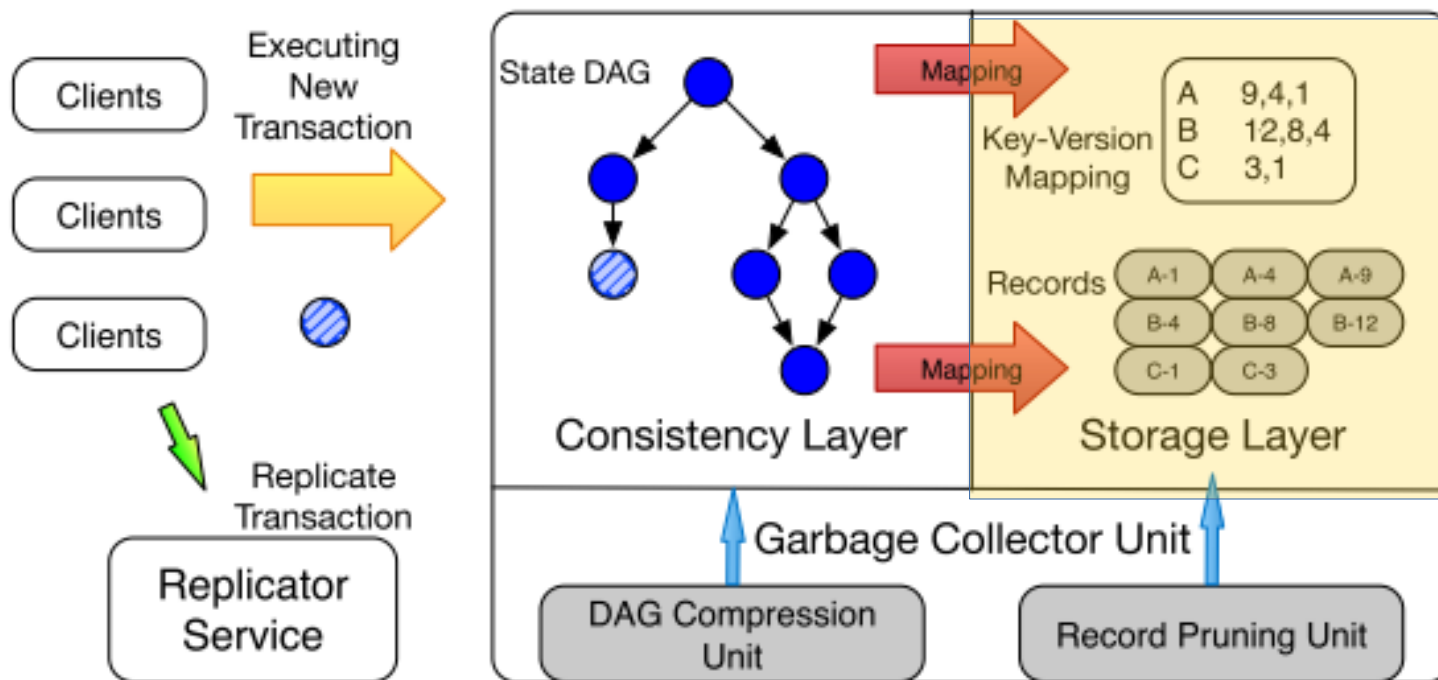
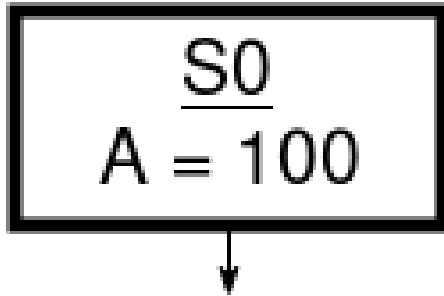


Figure 2: TARDiS architecture

Data structures

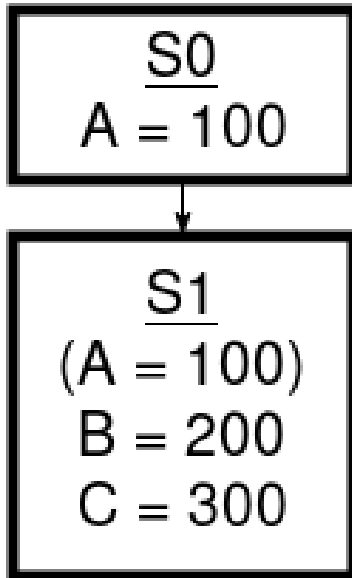


Key version mapping
A | S0

Record B-tree
A | S0

Fork paths:
The set of fork points
S0: {}

Data structures



Key version mapping

A | S0
B | S1
C | S1

Record B-tree

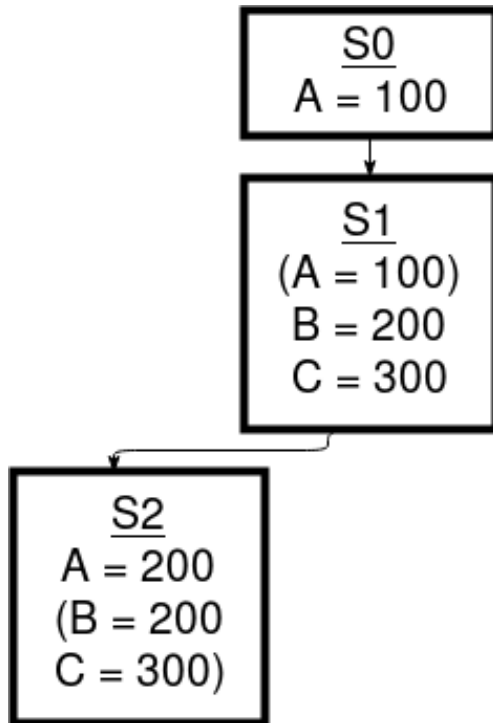
A | S0
B | S1
C | S1

Fork paths:

S0: {}

S1: {}

Data structures



Key version mapping

A | S2, S0

B | S1

C | S1

Record B-tree

A | S0 → S2

B | S1

C | S1

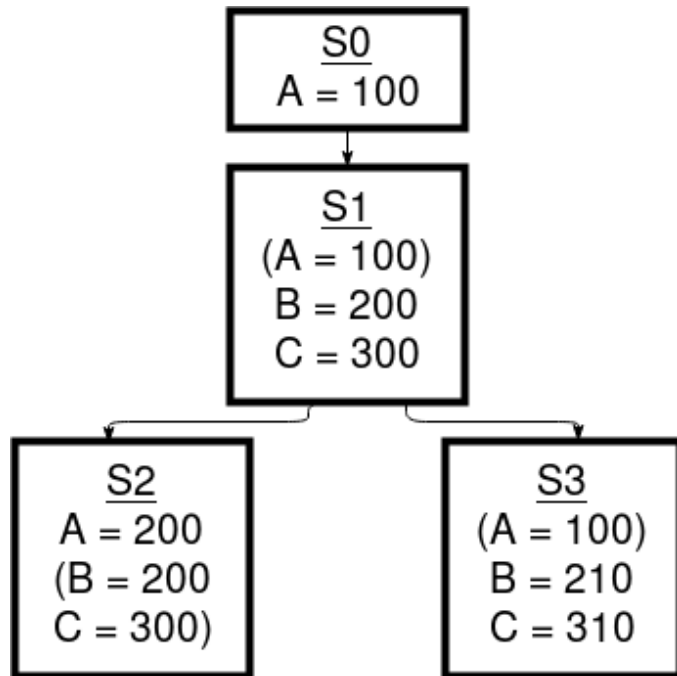
Fork paths: (set of tuples i, b where current state is both child of state i)

S0: {}

S1: {}

S2: { (1,1) }

Data structures



Key version mapping

A | S2, S0

B | S3, S1

C | S3, S1

Record B-tree

A | S0 → S2

B | S1 → S3

C | S1 → S3

Fork paths: (set of tuples i,b where current state is both child of state i)

S0,S1: { }

S2: { (1,1) }

S3: { (1,2) }

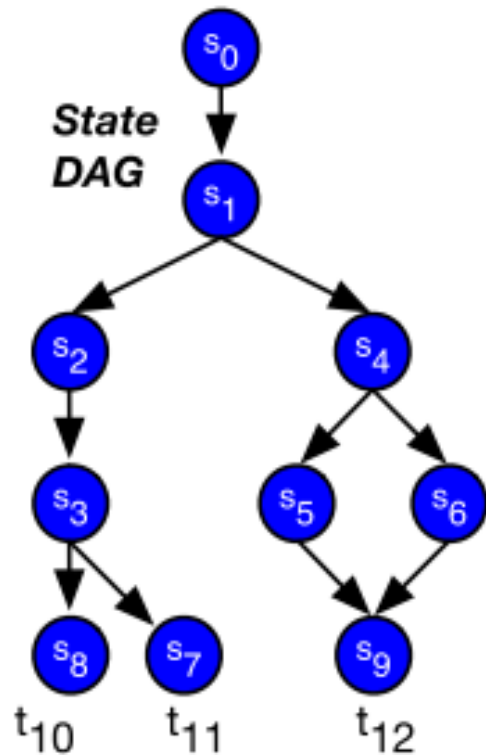
Data structures

```
1 descendantCheck(x, y):  
2   if x.id = y.id then return true  
3   else if x.id > y.id then return false  
4   else if x.path  $\not\subseteq$  y.path then return false  
5   else return true
```

Figure 7: Check if state y can see records associated with state x

A record version belongs to the selected branch if the fork path associated with this record version is a subset of the fork path of the transaction's read state

Data structures

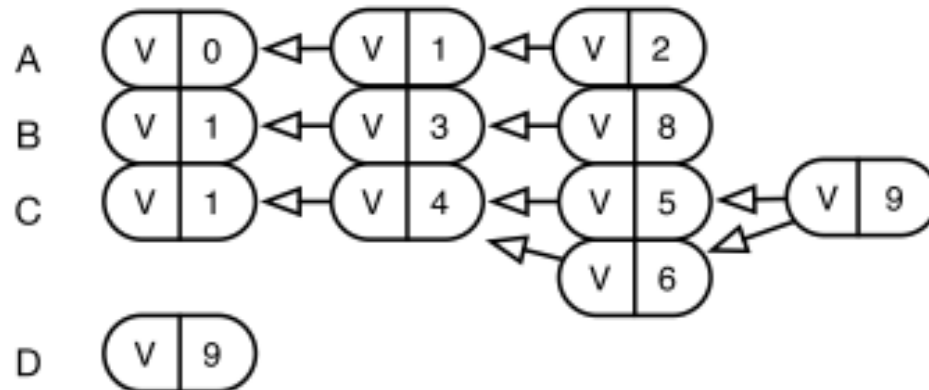


Key-Version Mapping

A	2,1,0
B	8,3,1
C	9,6,5,4,1
D	9

If transaction read state is S3
Then which record version of C?

Record B-Tree



TARDiS architecture

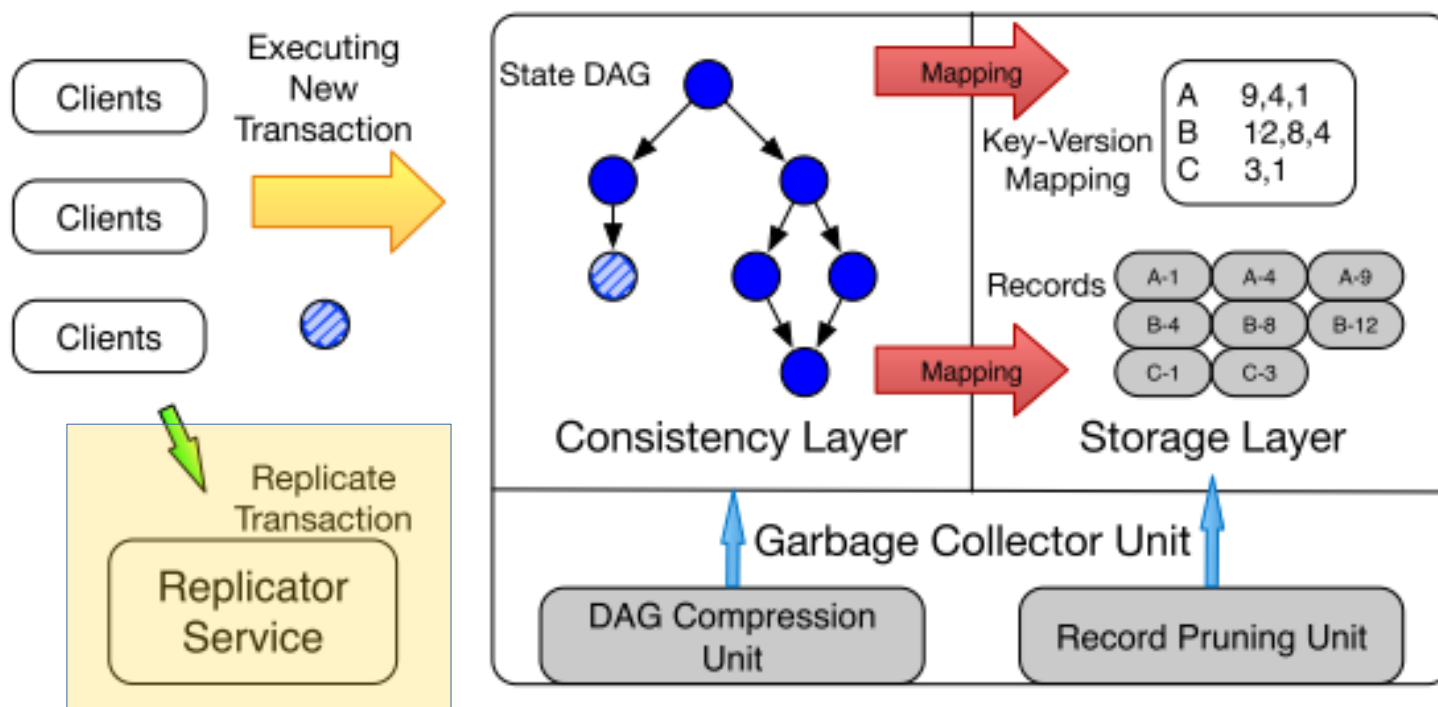


Figure 2: TARDiS architecture

Evaluation setup

Shared local cluster

2.67 GHz Intel Xeon CPU X5650

48GB memory

2Gbps network

3 dedicated server machines

3 dedicated replicators

Equally spread clients

For comparison

Databases

Berkley DB (BDB) - ACID datastore

An implementation that does not require read write transactions to be verified against read-only transactions (OCC)

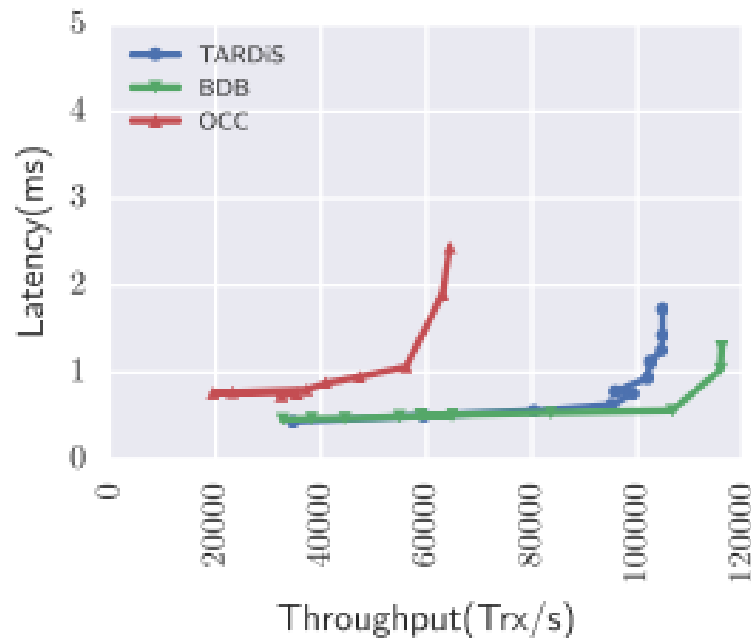
Operation composition

Read heavy (75R/25W)

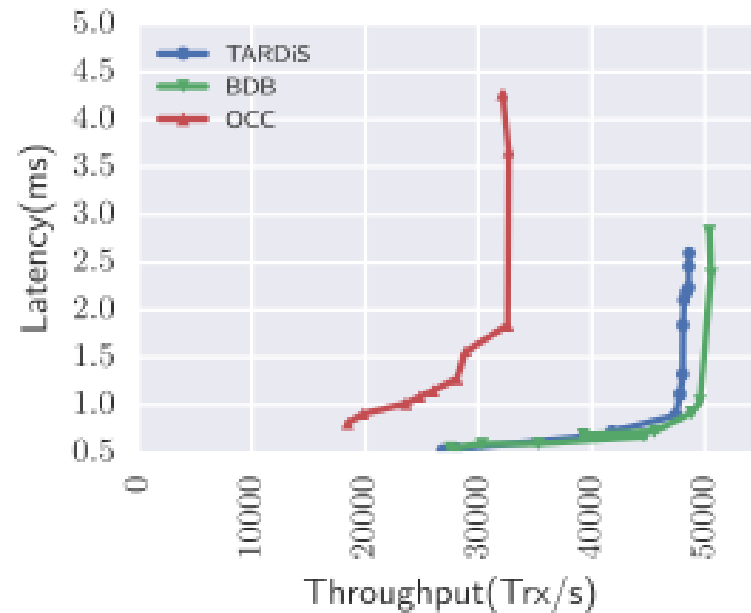
Write heavy (0R/100W)

Baseline TARDiS

Selecting constraints so that execution is serializable, and there is no branching



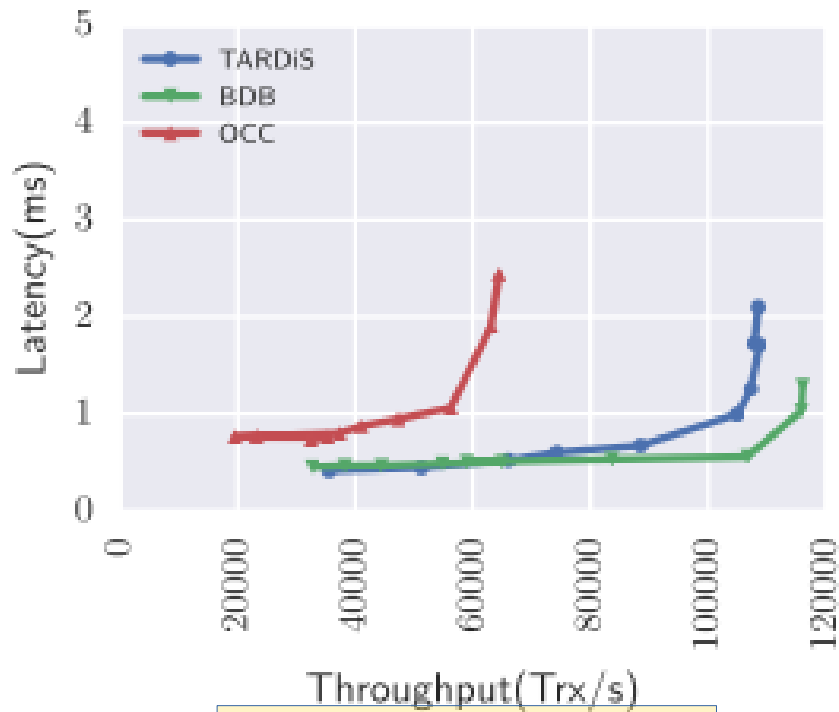
(a) Read-Heavy



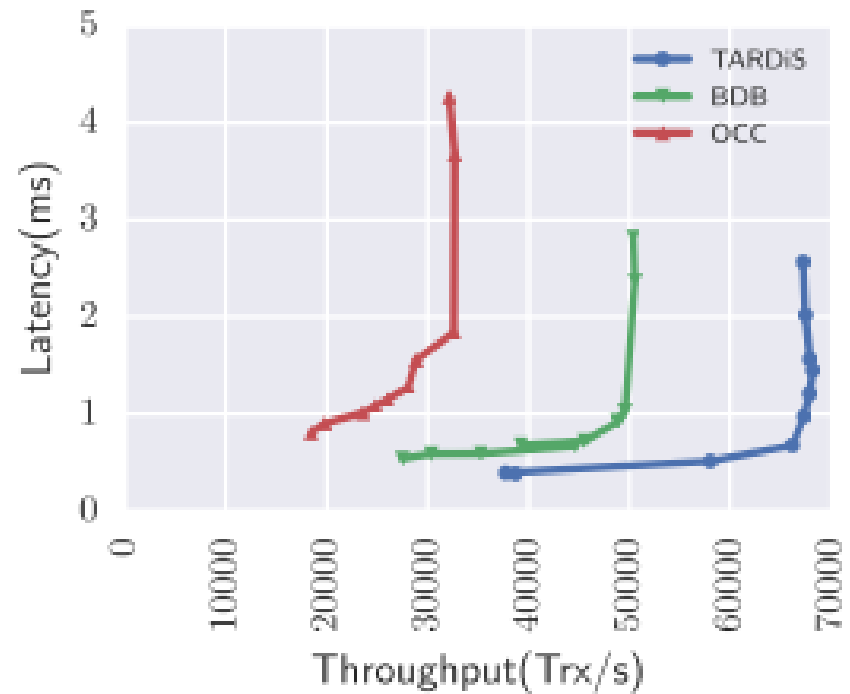
(b) Write-Heavy

Figure 9: TARDiS-BDB vs BerkeleyDB vs OCC

With branching

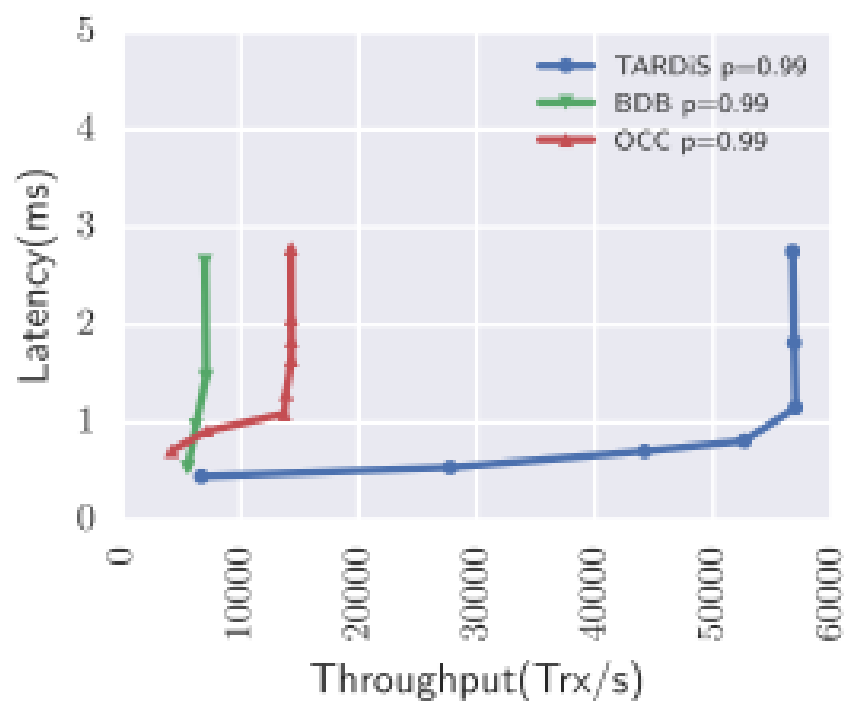


(a) Uniform Read-Heavy

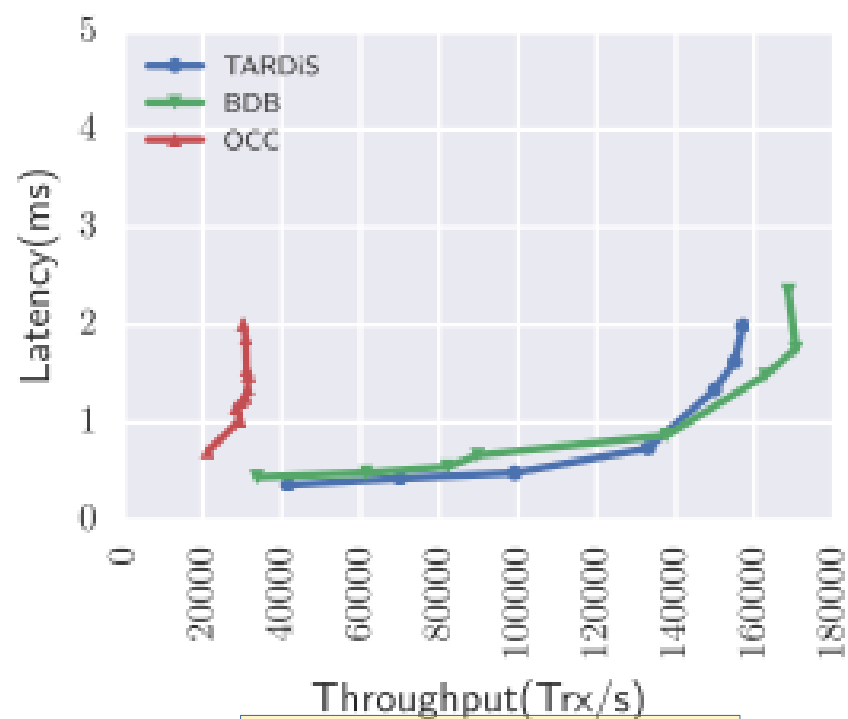


(b) Uniform Write-Heavy

With branching



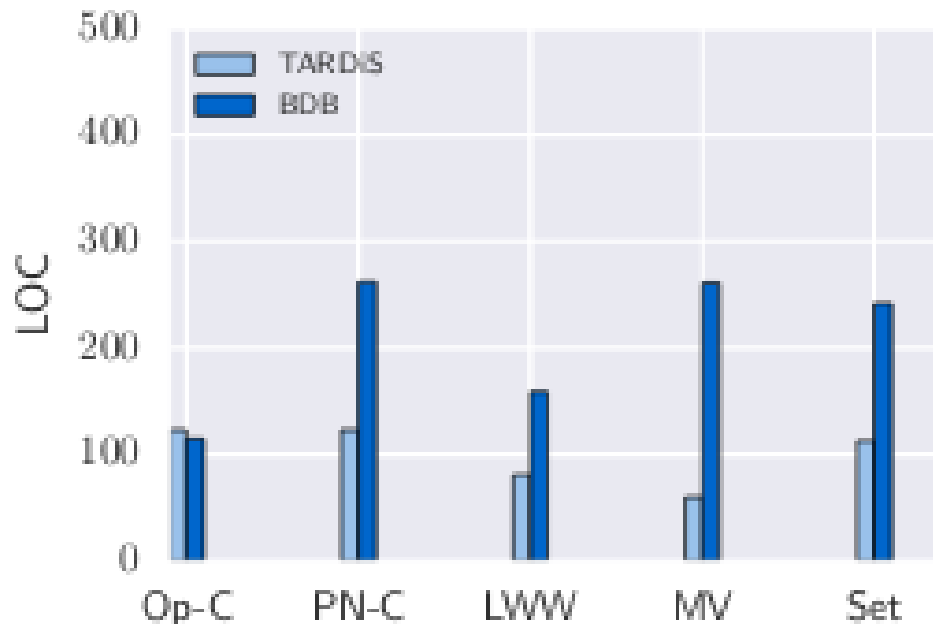
(c) Zipfian Write-Heavy



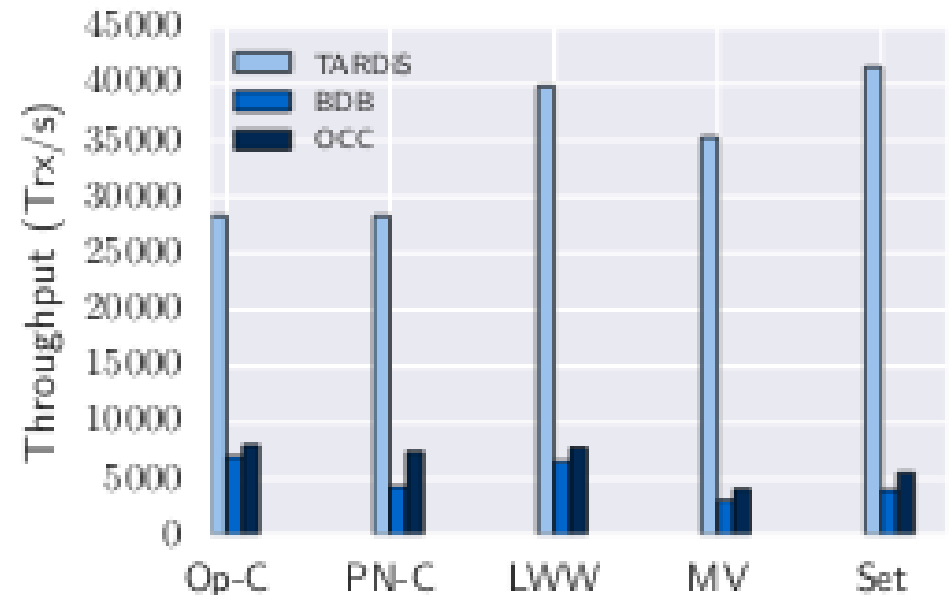
(d) Uniform Blind Writes

Figure 10: Benefit of branching as a function of workload

CRDT implementations



(a) CRDT Lines of Code



(b) CRDT Throughput

Op-C: Operation Based Counter, PN-C: State Based Counter, LWW: Last-Writer-Wins Register, MV: Multivalued Register, Set: Or-Set

Insight

Branching as a means to provide an abstraction that lifts WW conflicts to the application level so that application developer can determine the intended outcome of conflicts in a weakly consistent application

Next

Hard for programmer to reason about the whole application state in merge function. Therefore have the ability to compose a merge function from multiple merge functions

Having the ability to push and pull from other states so that synchronization can happen asynchronously and by on request