# Consistency in Transactional Distributed Databases: Protocols and Testing
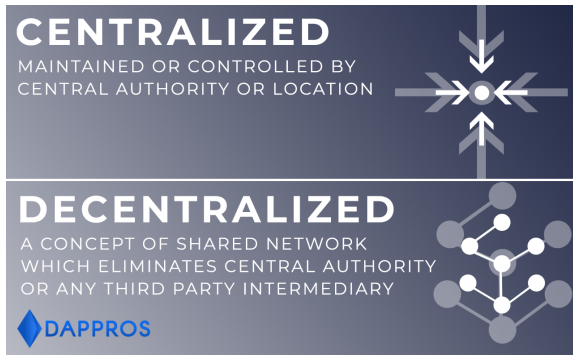
Hengfeng Wei (魏恒峰)

hfwei@nju.edu.cn
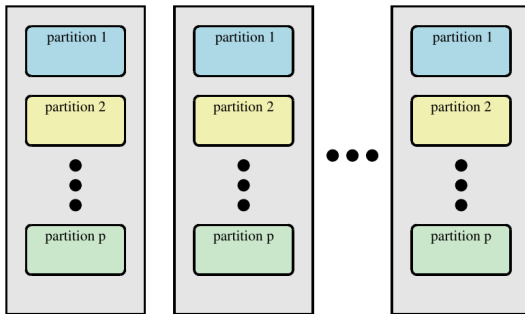
Nov. 04, 2022

Centralized Databases *vs.* Distributed Databases
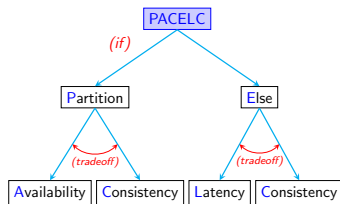
"Data Partition + Data Replication"



Data Consistency Problem

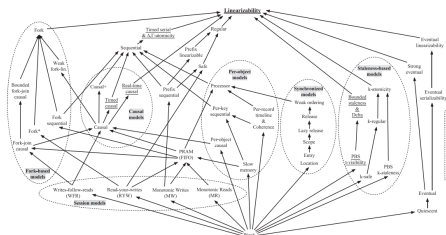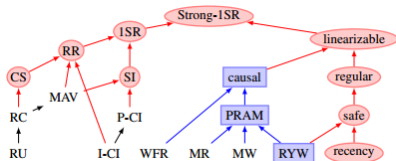(Strong) <u>C</u>onsistency, <u>A</u>vailability, <u>L</u>atency、 <u>P</u>atition tolerance



CAP Theorem



PACELC Tradeoff

Use consistency models to capture these tradeoffs

The theory of data consistency around consistency models:

Computability: What is possible or impossible?

Protocol: How to design fast, scalable, and
fault-tolerant protocols?

Testing: What is the complexity?

How to design efficient testing
algorihthms?

The theory of data consistency around consistency models:

Computability: What is possible or impossible?

Protocol: How to design fast, scalable, and fault-tolerant protocols?

Testing: What is the complexity?

How to design efficient testing algorihthms?

Classic problems with the ever-changing requirements

Research I (≥ 2012): Read/Write Register (读写寄存器)



分布式 NoSQL Key-Value 数据库 (TODO: 重新画图)
TODO: +research outcomes

Research II (≥ 2017): Replicated Data Types (复制数据类型)


(a) Google Docs


(b) Apache Wave

TODO: Jupiter/Redis/Riak


(c) Wikipedia


(d) LATEX Editor

TODO: +research outcomes

Research III ($\geq$ 2020): Distributed Transactions (分布式事务)
TODO: +research outcomes TODO: +logos

**UniStore: A fault-tolerant marriage of causal and strong consistency**

Manuel Bravo    Alexey Gotsman    Borja de Régil        Hengfeng Wei *
*IMDEA Software Institute*                    *Nanjing University*

ATC'2021 (CCF A)

UniStore is the first fault-tolerant and scalable transactional data store that combines causal and strong consistency.

*Partial Order-Restrictions Consistency (PoR consistency)*

TCC < PoR < SER

TCC: transactional causal consistency; SER: Serializability

Key Challenges (I): Ensure liveness in presence of faults

Key Challenges (II): Rigorous correctness proof

**UNISTORE: A fault-tolerant marriage of causal and strong consistency**

Manuel Bravo     Alexey Gotsman     Borja de Régil          Hengfeng Wei [*]

*IMDEA Software Institute*          *Nanjing University*

ATC'2021 (CCF A)

Fully responsible for the rigorous correctness proof:

▶ Finished a proof of 20 pages in the arXiv version

▶ Identified several nontrivial bugs in the early versions of the
  protocol[a]

---

[a]One of these bugs also exists in the well-known Granola protocol proposed by
James Cowling and Barbara Liskov, something that had gone unnoticed for 10
years.)

UNISTORE is a fast, scalable, and fault-tolerant

transactional distributed key-value store

that supports a combination of weak and strong consistency.

# What is UNISTORE?

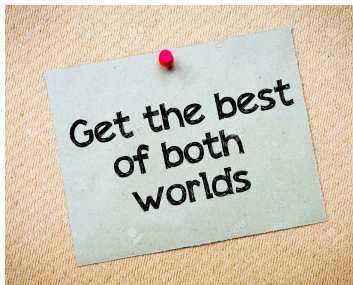UNISTORE is a fast, scalable, and fault-tolerant

transactional distributed key-value store

that supports a combination of weak and strong consistency.

Weak consistency: CAUSALCONSISTENCY

Strong consistency: SERIALIZABILITY

Weak consistency: low latency, high availability



Strong consistency: easy to preserve critical application invariants

# Why Uɴɪ-?

DEPOSIT   **WITHDRAW**   QUERY   INTEREST



Invariant: balance $\geq 0$

# Why Uni-?

DEPOSIT     WITHDRAW     QUERY     INTEREST



Invariant: balance $\geq 0$

Causal consistency allows two concurrent WITHDRAW
to execute without knowing each other.

# Why Uni-?

DEPOSIT  WITHDRAW  QUERY  INTEREST



Invariant: balance $\geq 0$

Causal consistency allows two concurrent WITHDRAW to execute without knowing each other.

Only WITHDRAW needs to use strong consistency.

UNISTORE implements a transactional variant of
Partial Order-Restrictions (POR) consistency [Li@ACT'2018]

(I)  transactional causal consistency by default
(II)  to specify conflicting transactions under strong consistency

# Consistency Model of UniStore

**Definition (Session Order)**

A transaction $t_1$ precedes a transaction $t_2$ in the session order, denoted $t_1 \xrightarrow{so} t_2$, if they are executed by the same client and $t_1$ is executed before $t_2$.

**Definition (Conflict Relation)**

The conflict relation, denoted $\bowtie$, between transactions is a symmetric relation.

$$t_1 \bowtie t_2 \iff t_2 \bowtie t_1.$$

# Consistency Model of UNISTORE

### Definition (PoR)

A set of transactions $T \triangleq T_{causal} \uplus T_{strong}$ committed by UNISTORE satisfies PoR if there exists a causal order $\prec$ on $T$ such that

# Consistency Model of UniStore

### Definition (PoR)

A set of transactions $T \triangleq T_{causal} \uplus T_{strong}$ committed by UniStore satisfies PoR if there exists a causal order $\prec$ on $T$ such that

$\qquad$ CAUSALITY: '$\prec$' is a partial order and $so \subseteq \prec$.

# Consistency Model of UNISTORE

Definition (PoR)

A set of transactions $T \triangleq T_{causal} \uplus T_{strong}$ committed by UNISTORE satisfies PoR if there exists a causal order $\prec$ on $T$ such that

CAUSALITY: '$\prec$' is a partial order and $so \subseteq \prec$.

CONFLICTORDERING: $\forall t_1, t_2 \in T_{strong}.\ t_1 \bowtie t_2 \implies t_1 \prec t_2 \lor t_2 \prec t_1$.

# Consistency Model of UNISTORE

**Definition (PoR)**

A set of transactions $T \triangleq T_{causal} \uplus T_{strong}$ committed by UNISTORE satisfies PoR if there exists a causal order $\prec$ on $T$ such that

CAUSALITY: '$\prec$' is a partial order and $so \subseteq \prec$.

CONFLICTORDERING: $\forall t_1, t_2 \in T_{strong}. \; t_1 \bowtie t_2 \implies t_1 \prec t_2 \lor t_2 \prec t_1$.

EVENTUALVISIBILITY: A transaction $t \in T$ that is either strong or originates at a correct data center eventually become visible at all correct data centers: from some point on, $t$ precedes in $\prec$ all transactions issued at correct data centers.

# Consistency Model of UNISTORE

## Definition (PoR)

A set of transactions $T \triangleq T_{causal} \uplus T_{strong}$ committed by UNISTORE
satisfies PoR if there exists a causal order $\prec$ on $T$ such that

CAUSALITY: '$\prec$' is a partial order and $so \subseteq \prec$.

CONFLICTORDERING: $\forall t_1, t_2 \in T_{strong}.\ t_1 \bowtie t_2 \implies t_1 \prec t_2 \vee t_2 \prec t_1$.

EVENTUALVISIBILITY: A transaction $t \in T$ that is either strong or
originates at a correct data center eventually
become visible at all correct data centers:
from some point on, $t$ precedes in $\prec$ all
transactions issued at correct data centers.

RETVAL: INTRETVAL $\wedge$ EXTRETVAL

# Consistency Model of UNISTORE

Consider a read $r$ from key $k$ in a transaction $t$.

INTRETVAL : read from the latest update on $k$ preceding $r$ in $t$

$$\text{RETVAL} = \text{INTRETVAL} \wedge \text{EXTRETVAL}$$

EXTRETVAL : read from the last update on $k$

of the latest transaction (in an order consistent with $\prec$) preceding $t$

# Consistency Model of UNISTORE

DEPOSIT  **WITHDRAW**  QUERY  INTEREST



Invariant: balance $\geq 0$

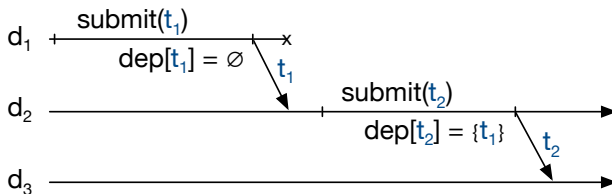Declaring that strong transactions
including WITHDRAW on the same account conflict.

To satisfy liveness (EventualVisibility) despite failures



A transaction $t \in T$ that is either strong or originates at a correct data center eventually become visible at all correct data centers.
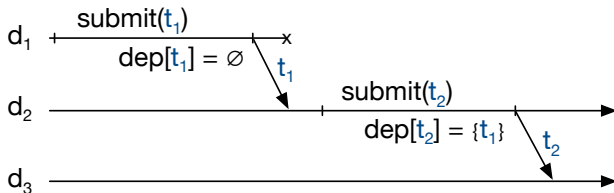
# Design Challenge of UNISTORE (I)

Data center $d_1$ crashes

before $t_1$ is replicated to correct data center $d_3$.
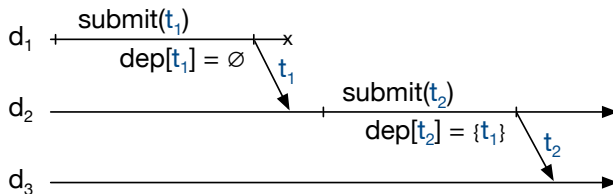
# Design Challenge of UNISTORE (I)

Data center $d_1$ crashes
before $t_1$ is replicated to correct data center $d_3$.



Transaction $t_2$ (at correct data center $d_2$)
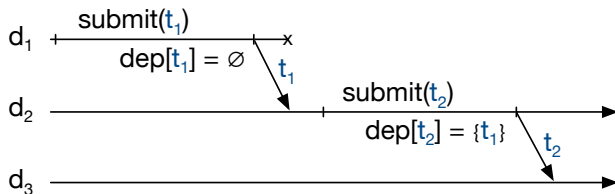may never become visible at correct data center $d_3$.

Data center $d_1$ crashes

before $t_1$ is replicated to correct data center $d_3$.
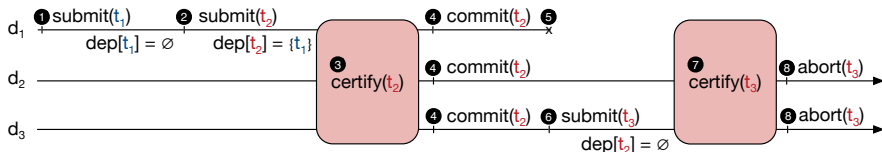
Data center $d_1$ crashes

before $t_1$ is replicated to correct data center $d_3$.



Data center $d_2$ need to forward causal transactions

to other data centers.

Data center $d_1$ crashes

before $t_1$ is replicated to correct data center $d_3$.

# Design Challenge of UNISTORE (II)
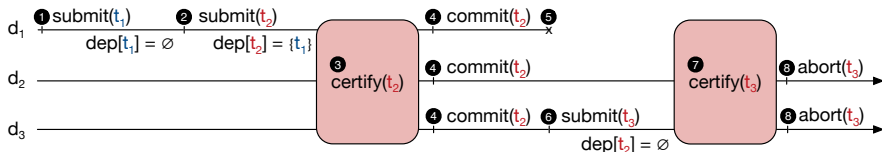
Data center $d_1$ crashes
before $t_1$ is replicated to correct data center $d_3$.



Transaction $t_2$ will never be visible at $d_3$.
No transaction $t_3$ conflicting with $t_2$ can commit
(by CONFLICTORDERING).

UniStore ensures that before a strong transaction commits,

all its causal dependencies are uniform,

i.e., will eventually become visible at all correct data centers.

UNISTORE ensures that before a strong transaction commits,

all its causal dependencies are uniform,

i.e., will eventually become visible at all correct data centers.



Transaction $t_1$ will eventually be visible at $d_3$.

Transaction $t_2$ will eventually be visible at $d_3$.

Transaction $t_3$ may be committed at $d_3$.

Causal transactions remain highly-available, i.e., committed locally.



A strong transaction may have to wait for some of its dependencies to become uniform before committing.

Causal transactions remain highly-available, i.e., committed locally.



A strong transaction may have to wait for some of its dependencies
to become uniform before committing.

However, this may cost too much.

# Performance of UNISTORE

UNISTORE makes a remote causal transaction visible to clients only after it is uniform.

Causal transactions are executed on an (almost) uniform snapshot that may be slightly in the past.

# Performance of UNISTORE

UNISTORE makes a remote causal transaction visible to clients
only after it is uniform.

Causal transactions are executed on an (almost)
uniform snapshot that may be slightly in the past.



A strong transaction only needs to wait for causal transactions
originating at the local data center to become uniform.

UNISTORE scales horizontally,
i.e., with the number of machines (partitions) in each data center.

# System Model

$$\mathcal{D} = \{1, \ldots, D\} : \text{the set of data centers}$$
$$\mathcal{P} = \{1, \ldots, N\} : \text{the set of (logical) partitions}$$



$p_d^m$ : the replica of partition $m$ at data center $d$

# System Model

$$D = 2f + 1 \text{ and } \leq f \text{ data centers may fail}$$



Any two replicas are connected by a reliable FIFO channel.

Messages between correct data centers will eventually be delivered.

# System Model

Replicas have loosely synchronized physical clocks.



The correctness of UNISTORE does not
depend on the precision of clock synchronization.

Fault-tolerant Causal Consistency Protocol

UNISTORE makes a remote causal transaction visible to clients only after it is uniform.

UNISTORE makes a remote causal transaction visible to clients only after it is uniform.

### Definition (Uniform)

A transaction is uniform if both the transaction and its causal dependencies are guaranteed to be eventually replicated at all correct data centers.

UNISTORE makes a remote causal transaction visible to clients
only after it is uniform.

## Definition (Uniform)

A transaction is uniform if both the transaction and its causal
dependencies are guaranteed to be eventually replicated at all
correct data centers.

A transaction is considered uniform
once it is visible at $f + 1$ data centers.

Each transaction is tagged with a commit vector *commitVec*.

$$commitVec \in [\mathcal{D} \to \mathbb{N}]$$

For a transaction originating at data center $d$,
we call $commitVec[d]$ its local timestamp.

Each transaction is tagged with a commit vector *commitVec*.

$$commitVec \in [\mathcal{D} \to \mathbb{N}]$$

For a transaction originating at data center $d$,
we call $commitVec[d]$ its local timestamp.

Commit vectors are sent to sibling replicas
via replication and forwarding.

Each replica $p_d^m$ maintains the following three vectors:

$$\mathsf{knownVec} \in [\mathcal{D} \to \mathbb{N}]$$

$$\mathsf{stableVec} \in [\mathcal{D} \to \mathbb{N}]$$

$$\mathsf{uniformVec} \in [\mathcal{D} \to \mathbb{N}]$$

# Metadata for Causal Transactions

$$\mathsf{knownVec} \in [\mathcal{D} \to \mathbb{N}]$$

Property (Property of knownVec)

For each data center $i$,

the replica $p_d^m$ stores the updates to partition $m$

by transactions originating at $i$ with local timestamps $\leq \mathsf{knownVec}[i]$.

# Metadata for Causal Transactions

$$\mathsf{stableVec} \in [\mathcal{D} \to \mathbb{N}]$$

**Property (Property of stableVec)**

For each data center $i$,

the data center $d$ stores the updates

by transactions originating at $i$ with local timestamps $\leq \mathsf{stableVec}[i]$.

1: **function** BROADCAST_VECS()
2:     **send** KNOWNVEC_LOCAL($m$, knownVec) **to** $p_d^l$, $l \in \mathscr{P}$
3:     **send** STABLEVEC($d$, stableVec) **to** $p_i^m$, $i \in \mathscr{D}$
4:     **send** KNOWNVEC_GLOBAL($d$, knownVec) **to** $p_i^m$, $i \in \mathscr{D}$

$$\text{stableVec} \in [\mathcal{D} \to \mathbb{N}]$$

5: **when received** KNOWNVEC_LOCAL($l, knownVec$)
6:     localMatrix[$l$] $\leftarrow knownVec$
7:     **for** $i \in \mathscr{D}$ **do**
8:         stableVec[$i$] $\leftarrow \min\{\text{localMatrix}[n][i] \mid n \in \mathscr{P}\}$
9:     stableVec[*strong*] $\leftarrow \min\{\text{localMatrix}[n][strong] \mid n \in \mathscr{P}\}$

# Metadata for Causal Transactions

$$\mathsf{uniformVec} \in [\mathcal{D} \to \mathbb{N}]$$

**Property (Property of uniformVec)**

All update transactions originating at $i$

with local timestamps $\leq \mathsf{uniformVec}[i]$

are replicated at $f + 1$ data centers including $d$.

# Metadata for Causal Transactions

1: **function** BROADCAST_VECS()
2:   **send** KNOWNVEC_LOCAL($m$, knownVec) **to** $p_d^l$, $l \in \mathscr{P}$
3:   **send** STABLEVEC($d$, stableVec) **to** $p_i^m$, $i \in \mathscr{D}$
4:   **send** KNOWNVEC_GLOBAL($d$, knownVec) **to** $p_i^m$, $i \in \mathscr{D}$

$$\text{uniformVec} \in [\mathcal{D} \to \mathbb{N}]$$

10: **when received** STABLEVEC($i$, $stableVec$)
11:   stableMatrix[$i$] $\leftarrow stableVec$
12:   $G \leftarrow$ all groups with $f + 1$ replicas that include $p_d^m$
13:   **for** $j \in \mathscr{D}$ **do**
14:     **var** $ts \leftarrow \max\{\min\{\text{stableMatrix}[h][j] \mid h \in g\} \mid g \in G\}$
15:     uniformVec[$j$] $\leftarrow \max\{\text{uniformVec}[j], ts\}$

# Metadata for Causal Transactions

$$\mathsf{uniformVec} \in [\mathcal{D} \to \mathbb{N}]$$

**Lemma**

*All update transactions*
*with commit vectors $\leq$ uniformVec are uniform.*

# Metadata for Causal Transactions

$$\mathsf{uniformVec} \in [\mathcal{D} \to \mathbb{N}]$$

**Lemma**

*All update transactions*
*with commit vectors $\leq$ uniformVec are uniform.*

UniStore makes a remote causal transaction visible to clients
only after it is uniform.

pastVec : causal past of client

1: **function** START()
2:      p ← a random partition in data center $d$
3:      $\langle \text{tid}, snapVec \rangle$ ← **send** START_TX(pastVec) **to** p
4:      pastVec ← $snapVec$
5:      **return** tid

$\forall i \in \mathcal{D} \setminus \{d\}$, all transactions originating at $i$
with local timestamps $\leq$ pastVec[$i$] are already uniform.

# Causal Consistency Protocol: Start

Causal transactions are executed on an (almost) <span style="color:red">uniform snapshot</span>.

1: **function** $\boxed{\text{START\_TX}(V)}$
2:     **for** $i \in \mathscr{D} \setminus \{d\}$ **do**
3:         $\boxed{\mathsf{uniformVec}[i] \leftarrow \max\{V[i], \mathsf{uniformVec}[i]\}}$
4:     **var** $tid \leftarrow \mathsf{generate\_tid}()$
5:     $\boxed{\mathsf{snapVec}[tid] \leftarrow \mathsf{uniformVec}}$
6:     $\mathsf{snapVec}[tid][d] \leftarrow \max\{V[d], \mathsf{uniformVec}[d]\}$
7:     $\mathsf{snapVec}[tid][strong] \leftarrow \max\{V[strong], \mathsf{stableVec}[strong]\}$
8:     **return** $\langle tid, \mathsf{snapVec}[tid] \rangle$

$\mathsf{snapVec}[tid][d]$ ensures "read-your-writes".

# Causal Consistency Protocol: Update

11: **function** UPDATE$(k, v)$
12:     _ $\leftarrow$ **send** DO_UPDATE(tid, $k, v$) **to** p
13:     **return** ok

17: **function** DO_UPDATE$(tid, k, v)$
18:     **var** $l \leftarrow$ partition$(k)$
19:     wbuff$[tid][l][k] \leftarrow v$
20:     rset$[tid][l] \leftarrow$ rset$[tid][l] \cup \{k\}$
21:     **return** ok

wbuff$[tid][l]$ : buffer for the latest local update on each key

# Causal Consistency Protocol: Read

6: **function** READ($k$)
7: $\quad \langle v, c \rangle \leftarrow$ **send** DO_READ(tid, $k$, lc) **to** p
8: $\quad$ **if** $c \neq \perp$ **then**
9: $\quad\quad$ lc $\leftarrow \max\{$lc, $c\}$
10: $\quad$ **return** $v$

9: **function** DO_READ($tid, k, c$)
10: $\quad$ lc $\leftarrow \max\{$lc, $c\}$
11: $\quad$ **var** $l \leftarrow$ partition($k$)
12: $\quad$ **if** wbuff$[tid][l][k] \neq \perp$ **then**
13: $\quad\quad$ **return** $\langle$wbuff$[tid][l][k], \perp\rangle$
14: $\quad \langle v, c \rangle \leftarrow$ **send** READ_KEY(snapVec$[tid], k$) **to** $p_d^l$
15: $\quad$ rset$[tid][l] \leftarrow$ rset$[tid][l] \cup \{k\}$
16: $\quad$ **return** $\langle v, c \rangle$

# Causal Consistency Protocol: Read

Causal transactions are executed on an (almost) uniform snapshot.

1: **when received** READ_KEY($snapVec, k$) **from** $p$
2:     **for** $i \in \mathscr{D} \setminus \{d\}$ **do**
3:         $\mathsf{uniformVec}[i] \leftarrow \max\{snapVec[i], \mathsf{uniformVec}[i]\}$
4:     **wait until** $\mathsf{knownVec}[d] \geq snapVec[d] \wedge \mathsf{knownVec}[strong] \geq snapVec[strong]$
5:     $\langle v, commitVec, c \rangle \leftarrow \mathsf{snapshot}(\mathsf{opLog}[k], snapVec)$   ▷ returns the latest $commitVec$ (in terms of Lamport clock order in Definition 50) such that $commitVec \leq snapVec$
6:     **send** $\langle v, c \rangle$ **to** $p$

**wait** : ensure that it is as up-to-date as required by the snapshot

# Causal Consistency Protocol: Commit

14: **function** COMMIT_CAUSAL_TX()
15:     $\langle vc, c \rangle \leftarrow$ **send** COMMIT_CAUSAL(tid, lc) **to** p
16:     pastVec $\leftarrow vc$
17:     lc $\leftarrow c$
18:     **return** ok

# Causal Consistency Protocol: Commit

Read-only transactions returns immediately.

22: **function** COMMIT_CAUSAL($tid, c$)
23:     $lc \leftarrow \max\{lc, c\} + 1$
24:     **if** $\forall l \in \mathscr{P}$. wbuff$[tid][l] = \emptyset$ **then**
25:         **return** $\langle$snapVec$[tid], lc\rangle$

26:     **var** $commitVec \leftarrow$ snapVec$[tid]$
27:     **send** PREPARE($tid$, wbuff$[tid][l]$, snapVec$[tid]$) **to** $p_d^l$, $l \in \mathscr{P}$
28:     **for all** $l \in \mathscr{P}$ **do**
29:         **wait receive** PREPARE_ACK($tid, ts$) **from** $p_d^l$
30:         $commitVec[d] \leftarrow \max\{commitVec[d], ts\}$
31:     **send** COMMIT($tid, commitVec, lc$) **to** $p_d^l$, $l \in \mathscr{P}$
32:     **return** $\langle commitVec, lc\rangle$

2PC protocol for update transactions

$ts$ : prepare timestamp from its local clock

7: **when received** PREPARE($tid, wbuff, snapVec$) **from** $p$
8:     **for** $i \in \mathscr{D} \setminus \{d\}$ **do**
9:         $\mathsf{uniformVec}[i] \leftarrow \max\{snapVec[i], \mathsf{uniformVec}[i]\}$
10:     **var** $ts \leftarrow \mathsf{clock}$
11:     $\mathsf{preparedCausal} \leftarrow \mathsf{preparedCausal} \cup \{\langle tid, wbuff, ts\rangle\}$
12:     **send** PREPARE_ACK($tid, ts$) **to** $p$

# Causal Consistency Protocol: Commit

**wait** : ensure that its local clock is up-to-date

13: **when received** COMMIT($tid, commitVec, c$)
14:     **wait until** clock $\geq commitVec[d]$
15:     $\langle tid, wbuff, \_ \rangle \leftarrow \mathsf{find}(tid, \mathsf{preparedCausal})$
16:     $\mathsf{preparedCausal} \leftarrow \mathsf{preparedCausal} \setminus \{\langle tid, \_, \_ \rangle\}$
17:     **for all** $\langle k, v \rangle \in wbuff$ **do**
18:         $\mathsf{opLog}[k] \leftarrow \mathsf{opLog}[k] \cdot \langle v, commitVec, c \rangle$
19:     $\mathsf{committedCausal}[d] \leftarrow \mathsf{committedCausal}[d] \cup \{\langle tid, wbuff, commitVec, c \rangle\}$

$\mathsf{committedCausal}[d]$ : for replication

# Causal Consistency Protocol: Replication

**Property (Property of knownVec)**

For each data center $i$,

the replica $p_d^m$ stores the updates to partition $m$

by transactions originating at $i$ with local timestamps $\leq$ knownVec[$i$].

```
1: function PROPAGATE_LOCAL_TXS()
2:     if preparedCausal = ∅ then
3:         knownVec[d] ← clock
4:     else
5:         knownVec[d] ← min{ts | ⟨_, _, ts⟩ ∈ preparedCausal} − 1
6:     var txs ← {⟨_, _, commitVec, c⟩ ∈ committedCausal[d] | commitVec[d] ≤ knownVec[d]}
7:     if txs ≠ ∅ then
8:         send REPLICATE(d, txs) to p_i^m, i ∈ 𝒟 \ {d}
9:         committedCausal[d] ← committedCausal[d] \ txs
10:    else
11:        send HEARTBEAT(d, knownVec[d]) to p_i^m, i ∈ 𝒟 \ {d}
```

HEARTBEAT : for liveness

# Adding Strong Transactions

# Requirement: CONFLICTORDERING

$$\forall t_1, t_2 \in T_{strong}. \ t_1 \bowtie t_2 \implies \ t_1 \prec t_2 \vee t_2 \prec t_1.$$

Each strong transaction is assigned a scalar strong timestamp.

$$commitVec \in [\mathcal{D} \cup \{strong\} \to \mathbb{N}]$$

# Metadata for Strong Transactions

$$\mathsf{knownVec} \in [\mathcal{D} \cup \{strong\} \to \mathbb{N}]$$

Property (Property of $\mathsf{knownVec}[strong]$)

Replica $p_d^m$ stores the updates to $m$ by all strong transactions with $commitVec[strong] \leq \mathsf{knownVec}[strong]$.

# Metadata for Strong Transactions

$$\text{stableVec} \in [\mathcal{D} \cup \{strong\} \to \mathbb{N}]$$

5: **when received** KNOWNVEC_LOCAL($l, knownVec$)
6:     localMatrix[$l$] $\leftarrow knownVec$
7:     **for** $i \in \mathcal{D}$ **do**
8:         stableVec[$i$] $\leftarrow \min\{\text{localMatrix}[n][i] \mid n \in \mathcal{P}\}$
9:     stableVec[$strong$] $\leftarrow \min\{\text{localMatrix}[n][strong] \mid n \in \mathcal{P}\}$

**Property (Property of stableVec[$strong$])**

Data center $d$ stores the updates by all strong transactions
with $commitVec[strong] \leq \text{knownVec}[strong]$.

# Metadata for Strong Transactions

$$\mathsf{uniformVec} \in [\mathcal{D} \to \mathbb{N}]$$

10: **when received** STABLEVEC($i, stableVec$)
11:     $\mathsf{stableMatrix}[i] \leftarrow stableVec$
12:     $G \leftarrow$ all groups with $f + 1$ replicas that include $p_d^m$
13:     **for** $j \in \mathcal{D}$ **do**
14:         **var** $ts \leftarrow \max\{\min\{\mathsf{stableMatrix}[h][j] \mid h \in g\} \mid g \in G\}$
15:         $\mathsf{uniformVec}[j] \leftarrow \max\{\mathsf{uniformVec}[j], ts\}$

The commit protocol for strong transactions
guarantees their uniformity.

# Strong Consistency Protocol: Commit

1: **function** COMMIT_STRONG($tid, c$)
2:    UNIFORM_BARRIER(snapVec[$tid$])
3:    $\langle d, vc, c \rangle \leftarrow$ CERTIFY($tid$, wbuff[$tid$], rset[$tid$], snapVec[$tid$], $c$)
4:    lc $\leftarrow \max\{$lc$, c\} + 1$
5:    **return** $\langle d, vc, $lc$\rangle$

A strong transaction only needs to wait for causal transactions
originating at the local data center to become uniform.

20: **function** UNIFORM_BARRIER($V, c$)
21:    lc $\leftarrow \max\{$lc$, c\} + 1$
22:    **wait until** uniformVec[$d$] $\geq V[d]$
23:    **return** lc

# Strong Consistency Protocol: Commit

1: **function** COMMIT_STRONG($tid, c$)
2:     UNIFORM_BARRIER($\mathsf{snapVec}[tid]$)
3:     $\langle d, vc, c \rangle \leftarrow$ CERTIFY($tid, \mathsf{wbuff}[tid], \mathsf{rset}[tid], \mathsf{snapVec}[tid], c$)
4:     $\mathsf{lc} \leftarrow \max\{\mathsf{lc}, c\} + 1$
5:     **return** $\langle d, vc, \mathsf{lc} \rangle$

$$\langle d \in \{\text{COMMIT}, \text{ABORT}\}, vc \rangle \leftarrow \text{CERTIFY}(t)$$

# Strong Consistency Protocol: Commit

```
1:  function COMMIT_STRONG(tid, c)
2:      UNIFORM_BARRIER(snapVec[tid])
3:      ⟨d, vc, c⟩ ← CERTIFY(tid, wbuff[tid], rset[tid], snapVec[tid], c)
4:      lc ← max{lc, c} + 1
5:      return ⟨d, vc, lc⟩
```

$$\langle d \in \{\text{COMMIT}, \text{ABORT}\}, vc \rangle \leftarrow \text{CERTIFY}(t)$$

**Multi-Shot Distributed Transaction Commit**

Gregory Chockler
Royal Holloway, University of London, UK

Alexey Gotsman[1]
IMDEA Software Institute, Madrid, Spain

### White-Box Atomic Multicast

Alexey Gotsman
IMDEA Software Institute

Anatole Lefort
Télécom SudParis

Gregory Chockler
Royal Holloway, University of London

2PC across partitions + Paxos among replicas of each partition

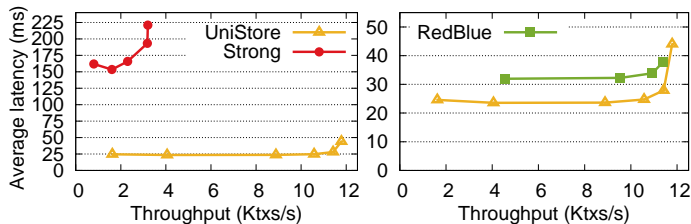uses white-box optimizations that minimize the commit latency

6: **upon** DELIVER_UPDATES($W$)
7:     **for** $\langle k, v, commitVec, c \rangle \in W$ in $commitVec[strong]$ order **do**
8:         $\mathsf{opLog}[k] \leftarrow \mathsf{opLog}[k] \cdot \langle v, commitVec, c \rangle$
9:         $\mathsf{knownVec}[strong] \leftarrow commitVec[strong]$

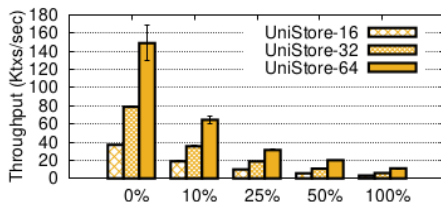# Evaluation

# Performance of UniStore

Throughput: 5% and 259% higher than RedBlue and Strong



RUBiS benchmark: throughput vs. average latency.

Latency: 24ms vs. 32ms of RedBlue and 162ms of Strong

# Scalability of UNISTORE



Scalability when varying the ratio of strong transactions.

UNISTORE is able to scales almost linearly.

For more evaluations, please refer to the paper.

# Conclusion

UNISTORE is a fast, scalable, and fault-tolerant

transactional distributed key-value store

that supports a combination of weak and strong consistency.

# Conclusion

UNISTORE is a fast, scalable, and fault-tolerant
transactional distributed key-value store
that supports a combination of weak and strong consistency.

"We expect the key ideas in UNISTORE to pave the way
for practical systems that combine causal and strong consistency."

# 总结

魏恒峰 (hfwei@nju.edu.cn)

| 聘期合同要求 | 工作情况 |
|---|---|
| **教学:** 承担一门课程 | 问题求解课程<br>五个学期; 共 164 学时<br>(2019 级本科生 "我心目中的好课程") |
| **科研:** 4-6 篇高水平论文 | 发表 3 篇 (含 1 篇短文)<br>在审 4 篇<br>(2017 年 CCF 优秀博士学位论文奖) |
| 人才培养 | 负责或协助指导学生 9 人次<br>(学术积累: 组织 TLA$^+$ 与 Coq 讨论班) |
| 主持/参与<br>多个基金项目 | 主持 1 项; 参与 1 项<br>个人可支配总经费 75 万元 |

Hengfeng Wei (hfwei@nju.edu.cn)