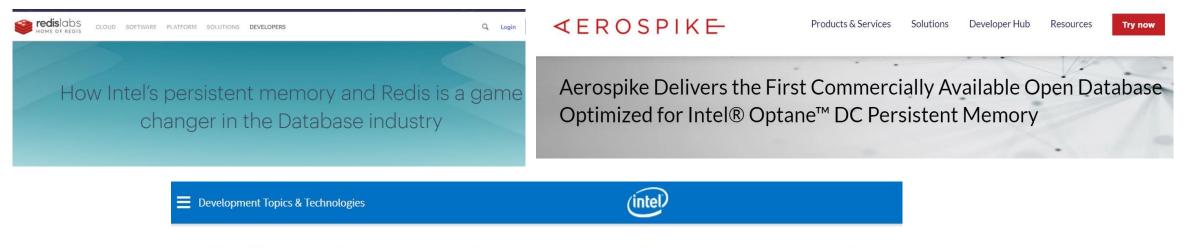# TIPS:Making Volatile Indexes Persistent With DRAM-NVMM Tiering

**R. Madhava Krishnan,**

Wook-hee Kim, Hee Won Lee[+*], Minsung Jang[†*],
Sumit Monga, Ajith Mathew, Changwoo Min

**VIRGINIA TECH**

+ *Samsung Electronics*

*† Peraton Labs*

# NVMM is Gaining Traction in Real-world Systems!

➢ Byte addressable Non-Volatile Main Memory (NVMM) has high capacity, low latency and durability

➢ Lots of interest in extending support for in-memory databases and key-value stores

intel OPTANE DC
PERSISTENT MEMORY

redislabs CLOUD SOFTWARE PLATFORM SOLUTIONS DEVELOPERS    Login

How Intel's persistent memory and Redis is a game changer in the Database industry

AEROSPIKE    Products & Services    Solutions    Developer Hub    Resources    Try now

Aerospike Delivers the First Commercially Available Open Database Optimized for Intel® Optane™ DC Persistent Memory

☰ Development Topics & Technologies    intel

Making NoSQL Databases Persistent-Memory-Aware: The Apache Cassandra* Example

# Porting Volatile Indexes for NVMM is Crucial!

➢ Index structures are core part of in-memory databases

➢ Recent research works focuses on converting volatile indexes to work on NVMM

➢ Manual porting is complex and error-prone

➢ Provides framework or guidelines to facilitate the porting

➢ State-of-the-art index conversion techniques

  ❏ NVTraverse [PLDI-20], PRONTO[ASPLOS-20], RECIPE[SOSP19]

# Existing Techniques Have a Narrow Scope

➢ Existing conversion techniques are proposed based on the concurrency control

- ❑ NVTraverse [PLDI-20] for lock-free indexes, e.g., Atomic CAS

- ❑ PRONTO [ASPLOS-20] for blocking indexes, e.g., Mutex

- ❑ RECIPE [SOSP-19] for fine-grained and lock-free indexes

## Existing Conversion Techniques Have Limited Applicability

4

# Existing Techniques Have Other Critical Limitations

➢ Support only Buffered Durable Linearizability [RECIPE]

➢ Not handling persistent memory leaks [RECIPE, NVTraverse]

➢ In-depth knowledge on the volatile index [RECIPE, NVTraverse]

➢ Can not scale beyond the DRAM capacity [PRONTO]

➢ High crash consistency overhead [PRONTO]

We propose TIPS to solve these problems and make the overall conversion process simple, intuitive and less error prone

# Talk Outline

➢ Motivation

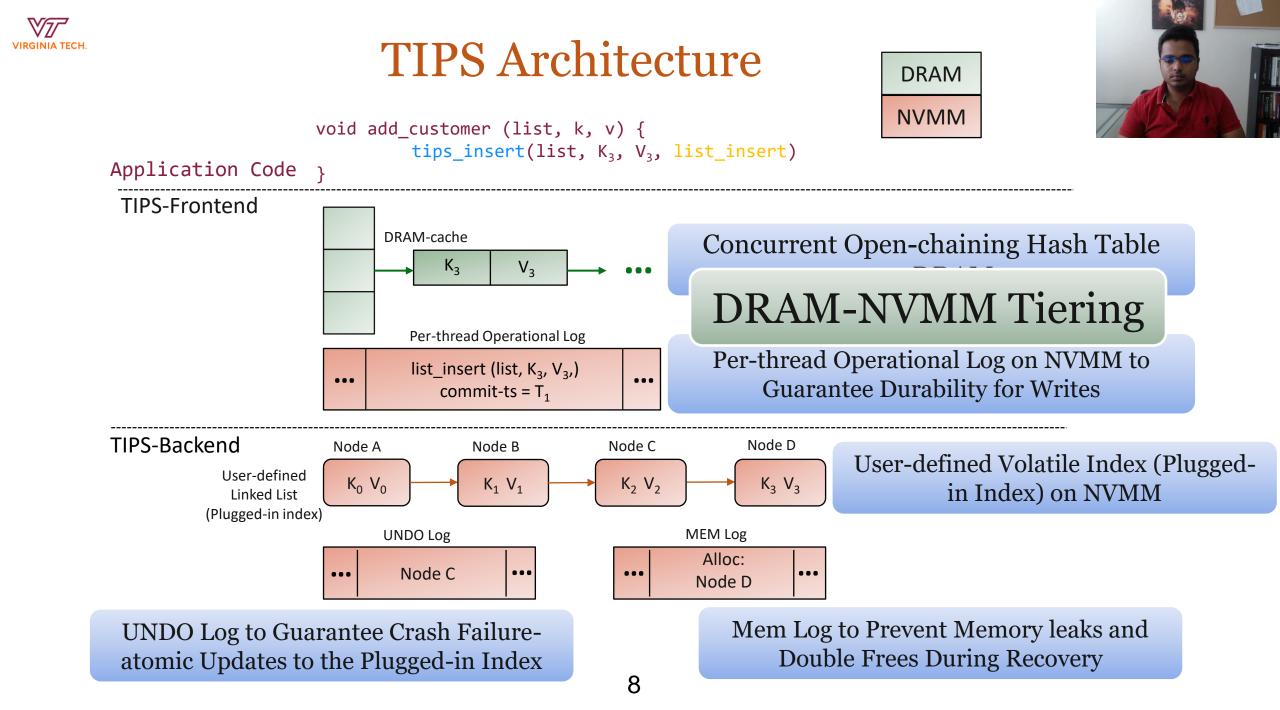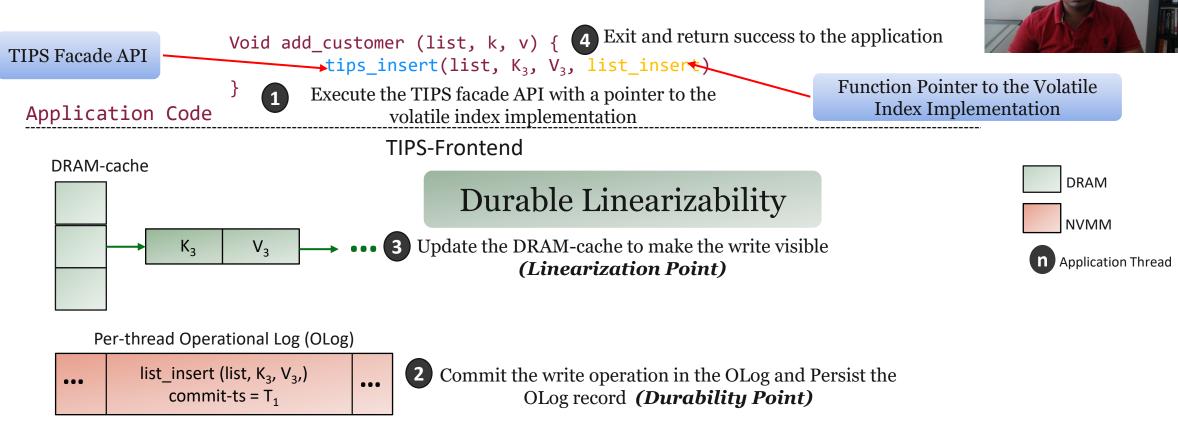➢ Overview

➢ Evaluation

➢ Conclusion

# Three Main Goals of TIPS

1) Support an Index-agnostic Conversion
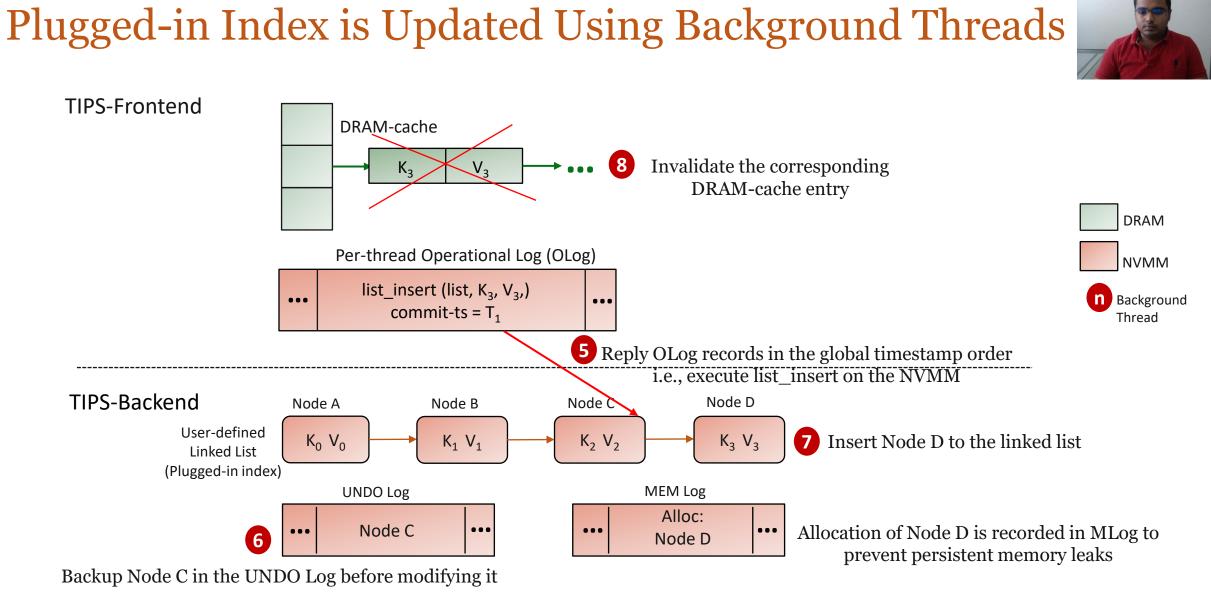
2) Guarantee Durable Linearizability for Correctness

3) Provide High-Performance and Scalability

# TIPS Architecture

```
void add_customer (list, k, v) {
        tips_insert(list, K_3, V_3, list_insert)
}
```

DRAM

NVMM

**Application Code**

**TIPS-Frontend**

DRAM-cache

$K_3$ | $V_3$ • • •

Concurrent Open-chaining Hash Table

DRAM-NVMM Tiering

Per-thread Operational Log

• • • | list_insert (list, $K_3$, $V_3$,) commit-ts = $T_1$ | • • •

Per-thread Operational Log on NVMM to Guarantee Durability for Writes

**TIPS-Backend**

Node A | Node B | Node C | Node D

User-defined Linked List (Plugged-in index)

$K_0$ $V_0$ → $K_1$ $V_1$ → $K_2$ $V_2$ → $K_3$ $V_3$

User-defined Volatile Index (Plugged-in Index) on NVMM

UNDO Log

• • • | Node C | • • •

MEM Log

• • • | Alloc: Node D | • • •

UNDO Log to Guarantee Crash Failure-atomic Updates to the Plugged-in Index

Mem Log to Prevent Memory leaks and Double Frees During Recovery

8

# Application Writes are Absorbed in TIPS-Frontend

```
Void add_customer (list, k, v) {        4  Exit and return success to the application
    tips_insert(list, K₃, V₃, list_insert)
}                                        1  Execute the TIPS facade API with a pointer to the
Application Code                           volatile index implementation
```

**TIPS Facade API**

**Function Pointer to the Volatile Index Implementation**

**TIPS-Frontend**

DRAM-cache

**Durable Linearizability**

$K_3$ | $V_3$  • • •  3  Update the DRAM-cache to make the write visible
*(Linearization Point)*

DRAM

NVMM

**n** Application Thread

Per-thread Operational Log (OLog)

| ... | list_insert (list, $K_3$, $V_3$,) commit-ts = $T_1$ | ... |

2  Commit the write operation in the OLog and Persist the OLog record  *(Durability Point)*

**Writes always happen at the fast TIPS-Frontend; Parallel disjoint writes regardless of concurrency model supported by the plugged-in index**

9

# Plugged-in Index is Updated Using Background Threads

**TIPS-Frontend**

DRAM-cache

$K_3$  $V_3$  • • •

**8** Invalidate the corresponding DRAM-cache entry

DRAM

NVMM

**n** Background Thread

Per-thread Operational Log (OLog)

• • • | list_insert (list, $K_3$, $V_3$,) commit-ts = $T_1$ | • • •

**5** Reply OLog records in the global timestamp order i.e., execute list_insert on the NVMM

**TIPS-Backend**

User-defined Linked List (Plugged-in index)

Node A | Node B | Node C | Node D
$K_0$ $V_0$ → $K_1$ $V_1$ → $K_2$ $V_2$ → $K_3$ $V_3$

**7** Insert Node D to the linked list

UNDO Log

**6** • • • | Node C | • • •

MEM Log

• • • | Alloc: Node D | • • •

Allocation of Node D is recorded in MLog to prevent persistent memory leaks

Backup Node C in the UNDO Log before modifying it

10

# Key Benefits of DRAM-NVMM Tiering

➢ Support index-agnostic conversion

  ❑ Allows plugged-in index to co-exist with the DRAM-cache

  ❑ No restrictions on the concurrency model of the volatile index

➢ Two different levels of concurrency (Tiered Concurrency Model)

  ❑ Concurrency model of DRAM-cache + Plugged-in index

  ❑ DRAM-Cache supports concurrent lock-free reads and disjoint writes

  ❑ Index with blocking concurrency (e.g., Mutex) can benefit from DRAM-cache

➢ Support Durable Linearizability agnostic of volatile index

# Can the TIPS-Backend Become a Scalability Bottleneck?

➢ TIPS-Frontend is fast and scalable with concurrent DRAM-cache and per-thread operational logging

➢ Backend writes are inherently slower because of

    ❑ Writes happening in the NVMM

    ❑ Notorious UNDO logging overhead

➢ Slower backend can easily bottleneck the frontend

➢ How do we make the TIPS-backend scalable?

# How TIPS Makes its Backend Scalable?

- ➢ A Key Intuition

    - ❑ Real-world workloads are rarely 100% writes

- ➢ We introduce two more techniques

UNO Logging Protocol to Reduce the UNDO Logging Overhead

Adaptive Scaling for Concurrent Background Writes

13

# UNO Logging Protocol

➢ All three logs (OLog, ULog, MLog) in TIPS works synergistically

➢ Not all modified addresses are required to be UNDO logged

    ❑ Selectively log only the addresses required for the correct recovery

➢ Perform UNDO logging only when the requested address

    ❑ is not previously UNDO-logged i.e., avoid redundant UNDO logging

    ❑ is not present in the OLog i.e., addresses that can not be recreated by OLog replay

➢ Significantly reduces the number of UNDO loggings performed

# Benefits of UNO Logging

➢ Makes the backend writes fast

   ❑ Number of UNDO logging is significantly reduced

   ❑ Enables write coalescing in the UNDO log

➢ Reduces crash consistency overhead in the write critical path

   ❑ Using OLog requires only 2 persist barriers

➢ Prevents persistent memory leaks

   ❑ Addresses in the MLog can be freed upon recovery

➢ UNO logging is index-agnostic

   ❑ applicable to any index irrespective of type or concurrency control

15

# Adaptive Scaling of Background Writers

- ➢ TIPS uses Adaptive Scaling to concurrently update the plugged-in index

    - ❏ Carefully orders the operations for a faster concurrent reply

- ➢ Adaptive scaling has some very nice properties

    - ❏ Automatically adjusts the worker count based on workload nature

    - ❏ Optimizes worker count based on the write-scalability

    - ❏ Prevents wastage of CPU cycles and other hardware resources

- ➢ Refer to the paper for more details and correctness

16

# Converting a Volatile Hash Table Using TIPS

```c
void hash_insert(hash_t *hash, key_t key, val_t value)
{
    node_t **pprev_next, *node, *new_node;
    int bucket_idx;
    pthread_rwlock_wrlock(&hash->lock);

    // Find a node in a collision list
    bucket_idx = get_bucket(key);
    node        = hash->buckets[bucket_idx]->head;
    pprev_next  = &hash->buckets[bucket_idx]->head;
    while (node && node->key < key) {
        pprev_next = &node->next;
        node = node->next;
    }
    // Case 1: update an existing key
    if (node->key == key) {
        // Before modifying the value, backup the old value
        node->value = value;
        tips_ulog_add(&node->value, sizeof(node->value))
    }   node->value = value; // then update then the node
        goto unlock_out;
    }
    // Case 2: add a new key
    // Allocate a new node using tips_alloc
    new_node = malloc(sizeof(*new_node));
    new_node = tips_alloc(sizeof(*new_node));
    new_node->key    = key;
    new_node->value = value;
    new_node->next   = node;
    // Before modifying the value, backup the old value
    *pprev_next = new_node;
    tips_ulog_add(pprev_next, sizeof(*pprev_next))
    *pprev_next = new_node; // then update then the node
unlock_out:
    pthread_rwlock_unlock(&hash->lock);
}
```

**M2**

**M1**

➤ Two simple guidelines for the conversion
- ❑ Replace the memory allocation/free with tips_alloc or tips_free
- ❑ Add tips_undo_add before modifying any NVMM address

➤ Key Benefits
- ❑ No need to manually insert flush/fence
- ❑ Makes the conversion simple and trivial
- ❑ Developers need not worry persistence and visibility ordering

17

# Other Interesting Designs

➢ Concurrency model and epoch-based GC in DRAM-cache

➢ Scan operation

➢ Adaptive Scaling

➢ UNO logging reclamation

➢ Recovery algorithm

➢ Detailed correctness section

# Talk Outline

- ➢ Motivation

- ➢ Overview

- ➢ Evaluation

- ➢ Conclusion

# Evaluation Questions

➢ How much LoC are required to convert an index using TIPS?

➢ How does TIPS perform against the prior index-specific conversion techniques?

➢ How does TIPS perform against the NVMM-optimized indexes?

# Evaluation Settings

➢ 2 socket server with Intel DCPMM

❑ 512GB NVMM and 64GB DRAM

❑ 2.4 GHZ 64 core Intel Xeon Gold CPU

➢ We evaluate 7 Indexes with different concurrency model

➢ YCSB with 32M keys for both integer and string type keys

| Workload Name | Read/Write/Scan Ratio | Workload Nature |
|---|---|---|
| Workload A | 50/50/0 | Write intensive |
| Workload B | 95/5/0 | Read intensive |
| Workload C | 100/0/0 | Read only |
| Workload D | 95/5/0 | Read Latest |
| Workload E | 0/5/95 | Short Range Scan |

# Evaluation Settings

➢ DRAM-cache size is set to 25% (300 MB)

➢ Compared against the state-of-the-art index conversion techniques

    ❑ PRONTO [ASPLOS-20]

    ❑ NVTraverse [PLDI-20]

    ❑ RECIPE [SOSP-19]

➢ And against NVMM-optimized indexes

    ❑ Hash Indexes- CCEH [FAST-19], LevelHashing [OSDI-18],

    ❑ B+Tree Indexes- FastFair [Fast-18 ], BzTree[VLDB-18]

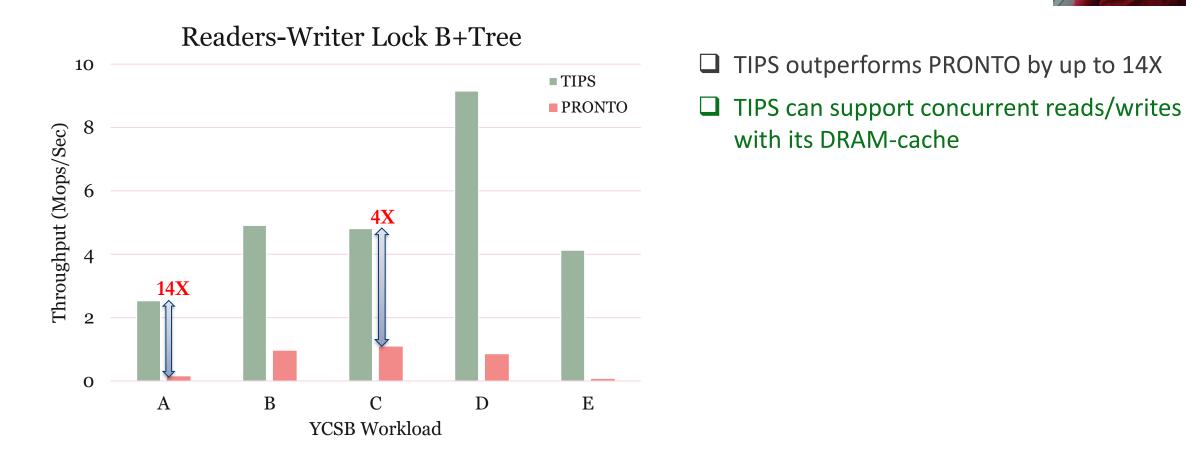    ❑ Radix Tree Indexes- WOART [FAST-17]

# LoC Required for Conversion

| Indexes | Concurrency Control | LoC change/ original LoC |
|---|---|---|
| Hash Table (HT) | Readers-Writer Lock | 5/211 |
| Lock-Free Hash Table (LFHT) | Non-blocking reads and writes | 5/199 |
| Binary Search Tree (BST) | Readers-Writer Lock | 5/203 |
| Lock-Free Binary Search Tree (LFBST) | Non-blocking reads and writes | 5/194 |
| B+Tree | Readers-Writer Lock | 8/711 |
| Adaptive Radix Tree (ART) | Non-blocking reads and blocking writes | 9/1.5k |
| Cache-Line Extensible Hash Table (CLHT) | Non-blocking reads and blocking writes | 8/2.8k |
| Redis Key-value Store | Blocking reads and writes | 18/10k |

## TIPS has better applicability and requires minimal code changes in the original codebase

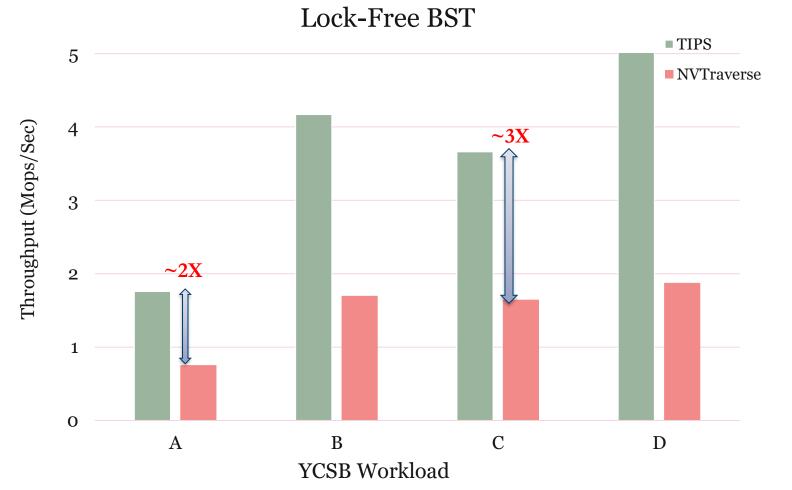# TIPS vs PRONTO for Blocking Indexes

## Readers-Writer Lock B+Tree



- ❑ TIPS outperforms PRONTO by up to 14X
- ❑ TIPS can support concurrent reads/writes with its DRAM-cache

Pronto: Easy and Fast Persistence for Volatile Data Structures [ASPLOS-2020]

# TIPS vs NVTraverse for Lock-Free Indexes

Lock-Free BST



❑ NVTraverse incurs 6 and 17 p-barriers for reads and writes

❑ TIPS incurs 2 p-barriers in the write critical path

❑ No p-barriers required for reads in TIPS

NVTraverse: In NVRAM Data Structures, the Destination Is More Important Than the Journey [PLDI-2020]

# Other Interesting Evaluations

➢ Performance comparison with the NVMM-optimized indexes

➢ Empirical analysis of TIPS design

➢ Scalability, skewness, large datasets etc.

➢ Sensitivity analysis

➢ Real-world application Redis

➢ More information on our conversion experience

# Conclusion

➢ Current Index conversion techniques

❑ Limited applicability

❑ Weak consistency guarantee

❑ Not address persistent memory leak

➢ **TIPS**

❑ No restrictions on concurrency model

❑ Offers strong consistency i.e., Durable Linearizability

❑ In addition to providing outstanding performance and scalability

❑ https://github.com/cosmoss-vt/tips

Thank You