# Tailwind: Fast and Atomic RDMA-based Replication

Yacine Taleb, Ryan Stutsman, Gabriel Antoniu, Toni Cortes



• General purpose in-memory key-value stores are widely used nowadays



 General purpose in-memory key-value stores are widely used nowadays

 Recent systems can process millions of requests/ second/machine: E.g. RAMCloud, FaRM, MICA, ...



• General purpose in-memory key-value stores are widely used nowadays

 Recent systems can process millions of requests/ second/machine: E.g. RAMCloud, FaRM, MICA, ...

 Key enablers: eliminating network overheads (e.g., kernel bypass) and leveraging multicore architectures



CPU is becoming a bottleneck in modern kv-stores
->Random memory access, key-value GET/PUT processing

CPU is becoming a bottleneck in modern kv-stores
->Random memory access, key-value GET/PUT processing

• Persistent in-memory kv-stores have to replicate data to survive failures



CPU is becoming a bottleneck in modern kv-stores
->Random memory access, key-value GET/PUT processing

• Persistent in-memory kv-stores have to replicate data to survive failures



CPU is becoming a bottleneck in modern kv-stores
->Random memory access, key-value GET/PUT processing

• Persistent in-memory kv-stores have to replicate data to survive failures



->Replication contends with normal request processing











• Replication impedes modern kv-stores

Replication impedes modern kv-stores and collocated applications



**Block Storage** 





**Graph Processing** 

**Stream Processing** 

• How to mitigate replication overhead?

- How to mitigate replication overhead?
- Techniques like Remote Direct Memory Access (RDMA) seem promising







## **Existing Protocols**

- Many systems use one-sided RDMA for replication
  - FaRM NSDI'14, DARE HPDC'15, HydraDB SC'16



# Outline

- Context
- Existing RDMA-based replication protocols
- Tailwind's design
- Evaluation
- Conclusion











## RDMA-based Replication Metadata

• A backup needs to insure the integrity of log-backup data + length of data



## RDMA-based Replication Metadata

The last checksum can cover all "previous" log entries

•



## RDMA-based Replication Metadata

The last checksum can cover all "previous" log entries

•













• The primary fails in the midst of updating metadata!





• Backup data unusable!



• The primary replicates A and corresponding metadata with a single RDMA



• The primary replicates B and corresponding metadata, right after A



• The primary fully replicates C, but partially the metadata



## RDMA-based Replication Fourth Try

• The primary replicates A and corresponding metadata



## RDMA-based Replication Fourth Try

• The primary partially replicates B, then fails



• The backup checks if objects were fully received





• The backup checks if objects were fully received



## Tailwind

- Keep the same client-facing interface (RPCs)
- Strongly-consistent primary-backup systems
- Appends only a 4-byte CRC32 checksum after each record
- Relies on **Reliable-Connected** queue pairs: messages are delivered at most once, in order, and without corruption
- Stop failures

• A primary chooses a backup, and requests an RDMA buffer



• A primary chooses a backup, and requests an RDMA buffer



 All subsequent replication requests are performed with one-sided RDMA WRITEs



• When the primary fills a buffer, it will chose a backup and repeat all steps



#### Tailwind: Failures

- Failures can happen at any moment
- RDMA complicates primary replica failures
- Secondary replica failures are naturally dealt with in storage systems













#### Failure scenarios: Partially Written Checksum

• Case 2: Partially transferred checksum





#### Failure scenarios: Partially Written Checksum

• Case 2: Partially transferred checksum





#### Failure scenarios: Partially Written Checksum

Case 2: Partially transferred checksum



#### Failure scenarios: Partially Written Object

• Case 3: Partially transferred object



# Outline

- Context
- Existing RDMA-based replication protocols
- Tailwind's design
- Evaluation
- Conclusion

#### Implementation

- Implemented on the RAMCloud in-memory kv-store
- Low latency, large scale, strongly consistent
- RPCs leverage fast networking and kernel bypass
- Keeps all data in memory, durable storage for backups

#### **RAMCloud Threading Architecture**



#### **Evaluation Configuration**

- Yahoo! Cloud Serving Benchmark (Workloads A (50%PUT), B(5%PUT), WRITE-ONLY)
- 20 million 100 bytes objects + 30 byte/key
- Requests generated with a Zipfian distribution
- RAMCloud replication vs Tailwind (<u>**3-way replication**</u>)

CPU	Xeon E5-2450 2.1 GHz 8 cores, 16 hw threads
RAM	16 GB 1600 MHz DDR3
NIC	Mellanox MX354A CX3 @ 56 Gbps
Switch	36 port Mellanox SX6036G
OS	Ubuntu 15.04, Linux 3.19.0-16, MLX4 3.4.0, libibverbs 1.2.1

#### **Evaluation Goals**

- How much CPU cycles can Tailwind save?
- Does it improve performance?
- Is there any overhead?

#### **Evaluation: CPU Savings**



#### **Evaluation: Latency Improvement**



#### Evaluation: Throughput



Even with 5% of PUTs, Tailwind increase throughput by 30%

#### Evaluation: Recovery Overhead



Recoveries with up to 10 million objects

#### Related Work

- One-sided RDMA systems: Pilaf ATC'13, HERD SIGCOMM'14, FaRM NSDI'14, DrTM SOSP'15, DrTM + R Eurosys'16, ...
- Mitigating replication overheads/Tuning consistency: RedBlue OSDI'12, Correctables OSDI'16
- Tailwind reduces replication CPU footprint and improves performance without sacrificing durability, availability, or consistency

## Conclusion

- Tailwind leverages one-sided RDMA to perform replication and leaves backups completely **idle**
- Provides backups with a protocol to protect against failure scenarios
- Reduces replication induced CPU usage while improving performance and latency
- Tailwind preserves client-facing RPC

# **Thank you! Questions?**