# Reducing Replication Bandwidth for Distributed Document Databases

**Lianghong Xu[1], Andy Pavlo[1], Sudipta Sengupta[2]
Jin Li[2], Greg Ganger[1]**
Carnegie Mellon University[1], Microsoft Research[2]

PARALLEL DATA LABORATORY
CARNEGIE MELLON UNIVERSITY

CARNEGIE MELLON
DATABASE GROUP

**Andy Pavlo**
@andy_pavlo

Today I am visiting @eliothorowitz at @mongodbinc to try to convince them to ditch MMAP & switch to anti-caching.



C.R.E.A.M.

CACHE RULES EVERYTHING AROUND ME

CARNEGIE MELLON
DATABASE GROUP

@andy_pavlo

RETWEETS
3

FAVORITES
3

9:57 AM - 3 Dec 2013

New York, NY

**#1** – You can sleep with grad students but <u>not</u> undergrads.

**#2** – Keep a bottle of water in your office in case a student breaks down crying.

**#3** – Kids <u>love</u> MongoDB, but they want to go work for Google.

# System Votes

| | |
|---|---|
| Google Spanner | 24 |
| mongoDB | 23 |
| redis | 10 |
| amazon DynamoDB | 5 |
| MySQL | 2 |
| APACHE HBASE | 1 |
| db Shards | 1 |

# Reducing Replication Bandwidth for Distributed Document Databases
*In ACM Symposium on Cloud Computing, pg. 1-12, August 2015.*

## More Info:

http://cmudb.io/doc-dbs

# Replication Bandwidth

# Replication Bandwidth

Primary
Database

**Goal: Reduce bandwidth for WAN geo-replication.**

Secondary

MMS

# Why Deduplication?

- Why not just **compress**?
  - *Oplog batches are small and not enough overlap.*

- Why not just use **diff**?
  - *Need application guidance to identify source.*

- **Deduplication** finds and removes redundancies.

# Compress vs. Dedup



Compression ratio

- Compress: 3.1
- trad-dedup: 9.1
- sDedup: 38.4
- sDedup + compress: 118.6

*20GB sampled Wikipedia dataset.*
*MongoDB v2.7 // 4MB Oplog batches*

# sDedup: Similarity Dedup



Client

Insertion & Updates

**Primary Node**

Oplog

Unsynchronized oplog entries

Source documents

**Delta Compressor**

Database

Source Document Cache

Deduplicated oplog entries

Oplog batch

**Secondary Node**

Oplog syncer

Source documents

**Delta Decompressor**

Database

Re-constructed oplog entries

Oplog

Replay

# Encoding Steps

- Identify Similar Documents
- Select the Best Match
- Delta Compression

# Identify Similar Documents

# Selecting the Best Match

## Initial Ranking

| Rank | Candidates | Score |
|------|------------|-------|
| 1 | Doc #2 | 2 |
| 1 | Doc #3 | 2 |
| 2 | Doc #1 | 1 |

## Final Ranking

| Rank | Candidates | Cached? | Score |
|------|------------|---------|-------|
| 1 | Doc #3 | Yes | 6 |
| 1 | Doc #1 | Yes | 3 |
| 2 | Doc #2 | No | 2 |

*Is doc cached?*

Source Document Cache

*If yes, reward 3x*

# Delta Compression

- Byte-level diff between source and target docs:
  – *Based on the xDelta algorithm*
  – *Improved speed with minimal loss of compression*

- **Encoding**:
  – *Descriptors about duplicate/unique regions + unique bytes*
- **Decoding**:
  – *Use source doc + encoded output*
  – *Concatenate byte regions in order*

# Evaluation

- MongoDB setup (v2.7)
  - *1 primary, 1 secondary node, 1 client*
  - *Node Config: 4 cores, 8GB RAM, 100GB HDD storage*

- Datasets:
  - *Wikipedia dump (20GB out of ~12TB)*
  - *Stack Exchange data dump (10GB out of ~100GB)*

# Compression: Wikipedia

■ sDedup    ■ trad-dedup



*20GB sampled Wikipedia dataset*

# Memory: Wikipedia

■ sDedup    ■ trad-dedup



Memory (MB)

| Chunk Size | sDedup | trad-dedup |
|---|---|---|
| 4KB | 34.1 | 80.2 |
| 1KB | 47.9 | 133.0 |
| 256B | 57.3 | 272.5 |
| 64B | 61.0 | 780.5 |

*20GB sampled Wikipedia dataset*

# Compression: StackExchange

■ sDedup    ■ trad-dedup



*10GB sampled StackExchange dataset*

# Throughput Overhead



Legend: with sDedup (pink), w/o sDedup (blue)

Left bar chart — Insertion throughput (MB/s):
- Wikipedia: 15.1 (with sDedup), 15.2 (w/o sDedup)
- Stack Exchange: 9.2 (with sDedup), 9.3 (w/o sDedup)

Right charts — Insertion throughput (MB/s) vs Run time (seconds): Wikipedia, Stack Exchange

# Dedup + Sharding



*20GB sampled Wikipedia dataset*

# Failure Recovery



Failure Point

Normal
Failure

Inserted documents (thousand)

Compression ratio

*20GB sampled Wikipedia dataset.*

# Conclusion

- Similarity-based deduplication for replicated document databases.

- **sDedup** for MongoDB (v2.7)
  - *Much greater data reduction than traditional dedup*
  - *Up to 38x compression ratio for Wikipedia*
  - *Resource-efficient design for inline deduplication with negligible performance overhead*

# What's Next?

- Port code to MongoDB v3.1
- Integrating **sDedup** into WiredTiger storage manager.
- Need to test with more workloads.

- Try not to get anyone pregnant.

# WiredTiger vs. sDedup

|  | Compression Ratio |
|---|---|
| Snappy | 1.6x |
| zLib | 3.0x |
| sDedup (no compress) | 38.4x |
| sDedup + Snappy | 60.8x |
| sDedup + zLib | 114.5x |

*20GB sampled Wikipedia dataset.*

# END

@andy_pavlo