# Constructing and Analyzing the LSM Compaction Design Space

BOSTON UNIVERSITY

Subhadeep Sarkar    Dimitris Staratzis    Zichen Zhu    Manos Athanassoulis

## LSM Compactions

### LSMs everywhere

RocksDB  WT  cassandra  accumulo™
QuasarDB  influxdb  SQLite

Fast Ingestion
+
Good Reads
+
Great Space Utilization

▶ **Periodic Compaction**

### Full Compaction vs. Partial Compaction

compact all data in a level          compact only part of a level



❶ insert 18   ❷ flush buffer & compact with Level 1
❸ Level 1 saturated; initiate compaction
❹ compact all data in Level 1 with Level 2

❶ insert 18   ❷ flush buffer & compact with overlapping files
❸ Level 1 saturated; initiate compaction
❹ compact one file from Level 1 with overlapping files in Level 2

⎯ memory buffer   ⌐ level capacity   ☐ file   ☐ file to compact   ☐ files after compaction

## Design Choice

### The Compaction Black Box

workload + LSM-tuning → ◼ → performance

**Which compaction strategy to use ❓**

### The Design Questions ▶ Design Primitives

- **How** to lay out the data? ▶ **Data Layout**
- **How much** data to move? ▶ **Compaction Granularity**
- **Which** data block to be moved? ▶ **Data Movement Policy**
- **When** to re-organize layout? ▶ **Compaction Trigger**

## The Design Space

### Data Layout (D)

leveling
tiering
hybrid

### Granularity (G)

levels
file(s)
runs

### Data Movement Policy (D)

file(s)

**Which file(s) to pick?**
- Round-robin
- Least overlap
- Coldest
- Oldest
- Most tombstones

### Trigger (T)

saturation
#runs

- Space ampl.
- Age of a file
- Read ampl.

### Primitive Ensemble

**Full**
- Ⓛ leveling
- Ⓖ levels
- Ⓓ N/A
- Ⓣ level saturation

**1-Lvl**
- Ⓛ hybrid leveling
- Ⓖ file
- Ⓓ least overlap w/ parent
- Ⓣ level saturation

**Tier**
- Ⓛ tiering
- Ⓖ runs (in same level)
- Ⓓ N/A
- Ⓣ #runs / space ampl.

## Observation

### Experimental Results



Data moved (GB) vs Compaction strategies — 45%, Raw data size; strategies: Full, LO+1, LO+2, RR, Cold, Old, Tier, 1–Lvl. Legend: read, write



Tail write latency (ms) vs Compaction strategies — 5x; strategies: Full, LO+1, LO+2, RR, Cold, Old, Tier, 1–Lvl.

- Smaller compaction granularity reduces data movement
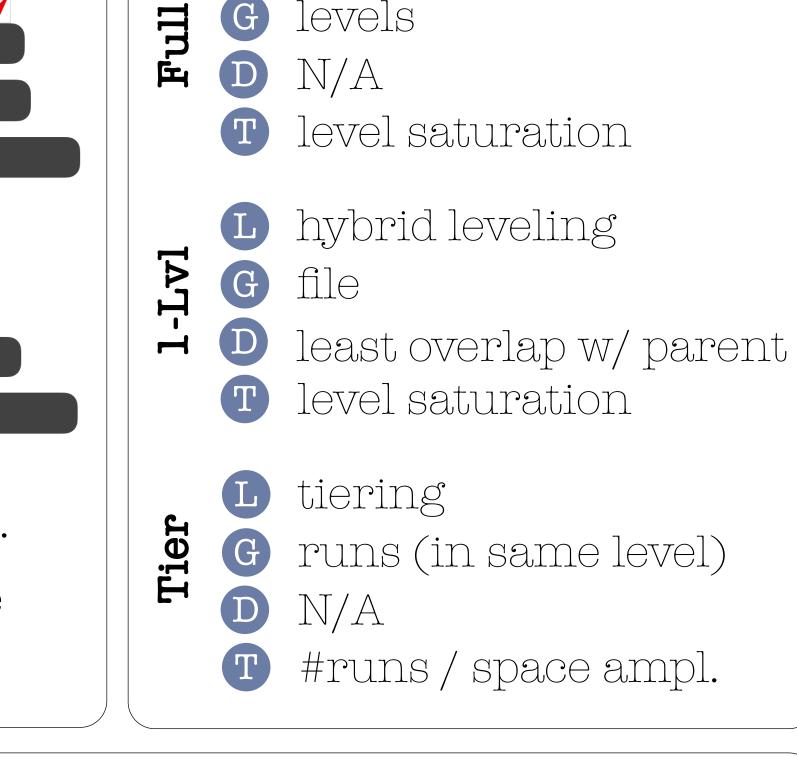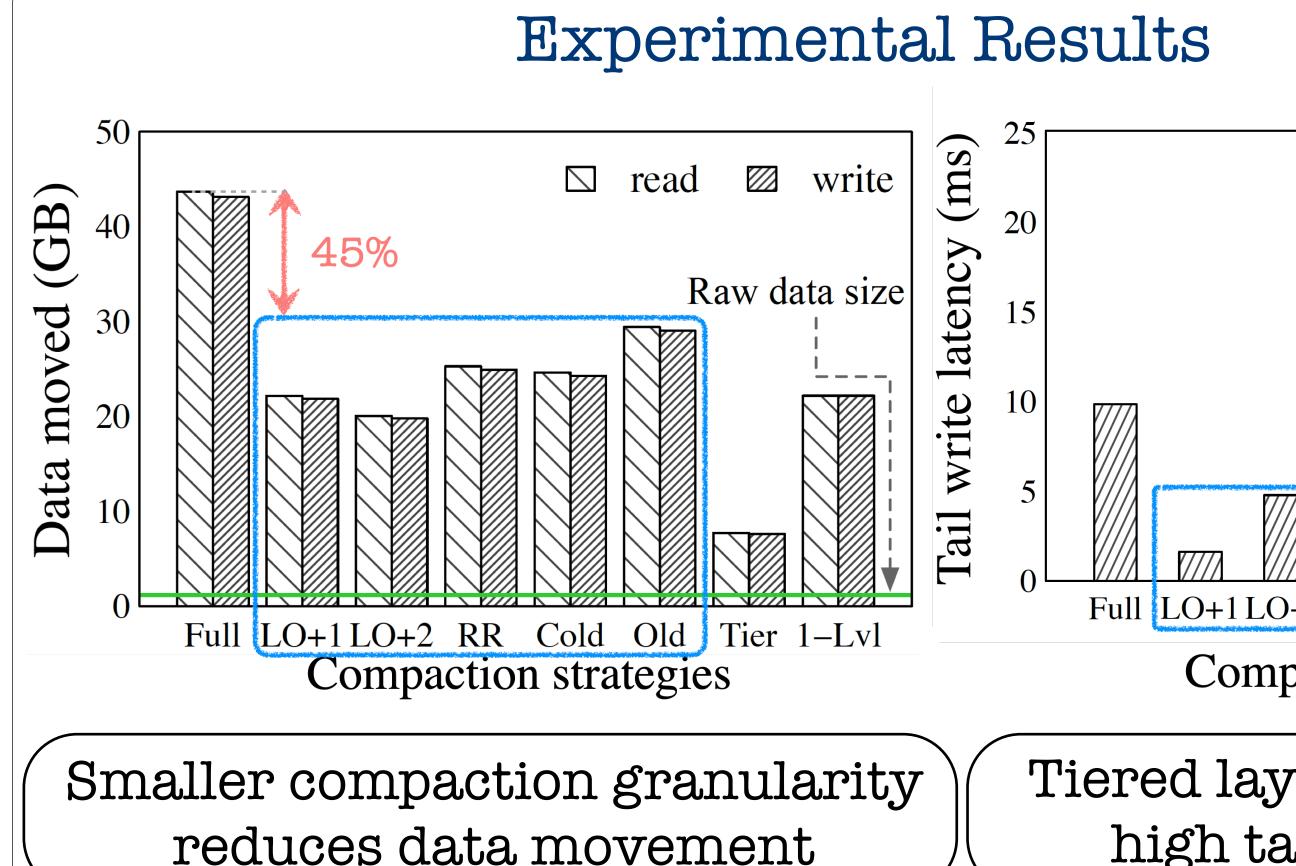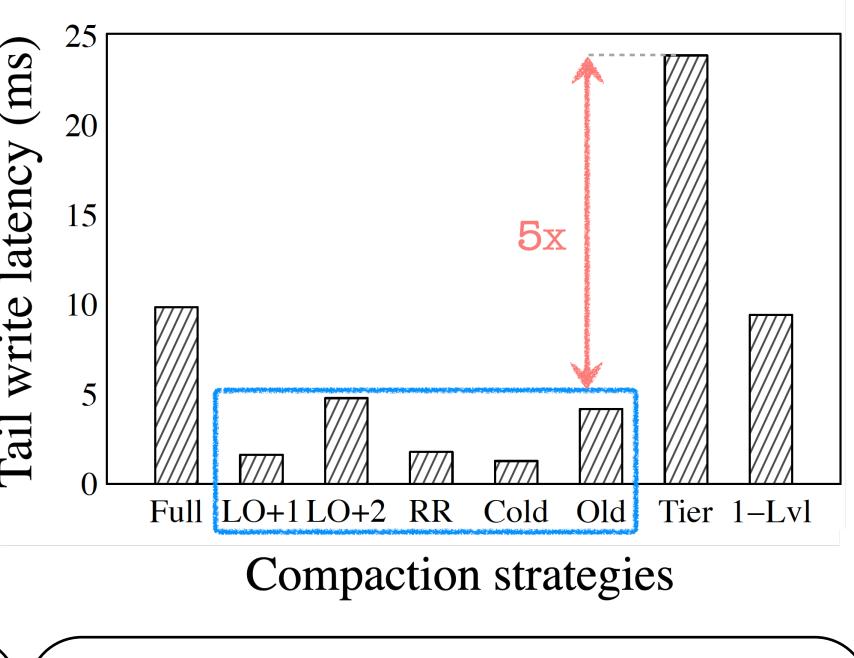- Tiered layout has extremely high tail write latency

### Key Takeaways

- Know your LSM-compaction
- Avoid the worst designs
- Adapt choices w/ workloads
- Design new compactions