

Towards Accurate and Fast Evaluation of Multi-Stage Log-Structured Designs

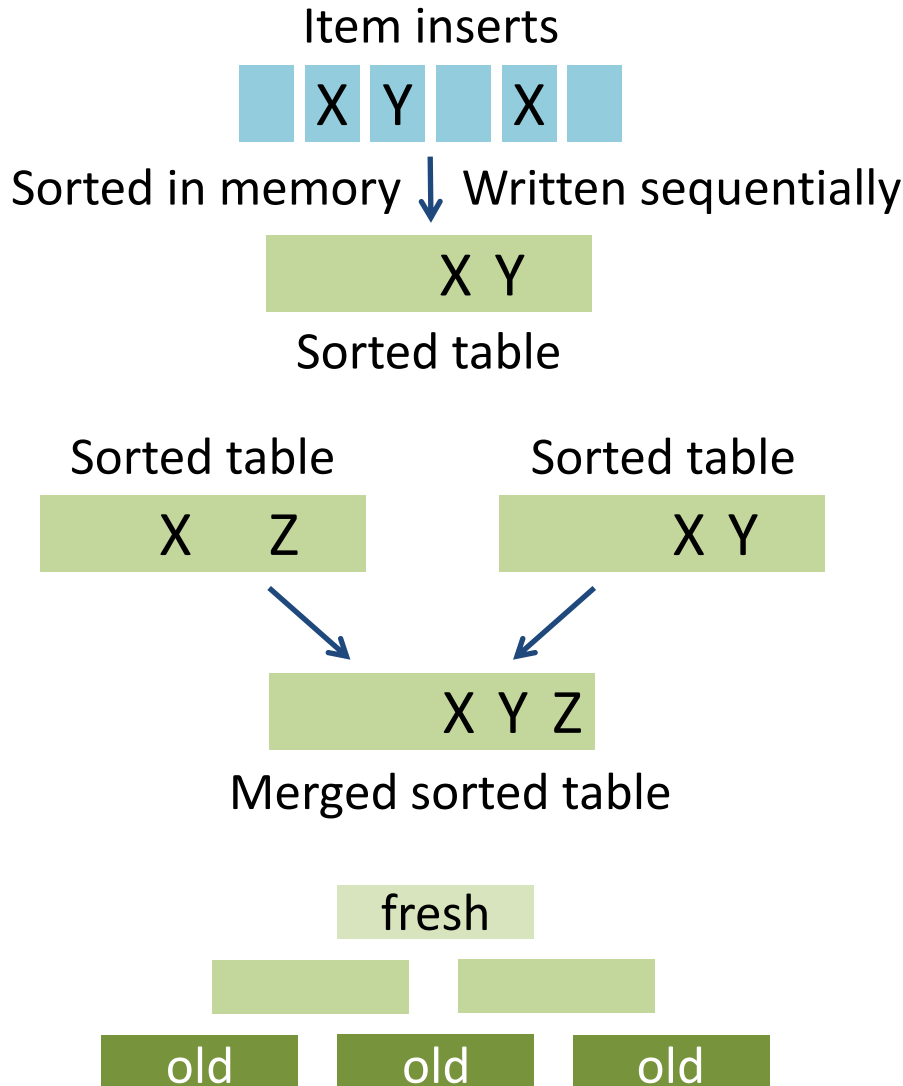
Hyeontaek Lim

David G. Andersen, Michael Kaminsky[†]

Carnegie Mellon University

[†]Intel Labs

Multi-Stage Log-Structured (“MSLS”) Designs



Example: LevelDB, RocksDB, Cassandra, HBase, ...

(Naïve) Log-structured design

- ⇒ Fast writes with sequential I/O
- ⇒ Slow query speed
- ⇒ Large space use

Compaction

- ⇒ Fewer table count
- ⇒ Less space use
- ⇒ Heavy I/O required

Multi-stage design

- ⇒ Cheaper compaction
by segregating fresh and old data

MSLS Design Evaluation Needed



Mobile app



Filesystem



Desktop app

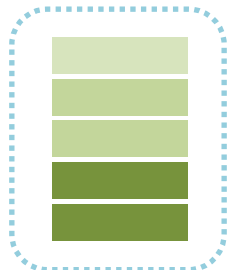


Data-intensive computing

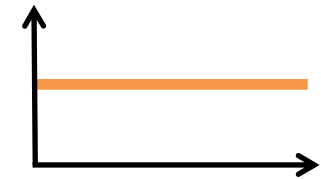
Problem: How to evaluate and tune MSLS designs for a workload?



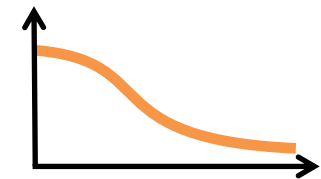
Large design space



Many tunable knobs



Diverse workloads



Two Extremes of Prior MSLS Evaluation

Speed

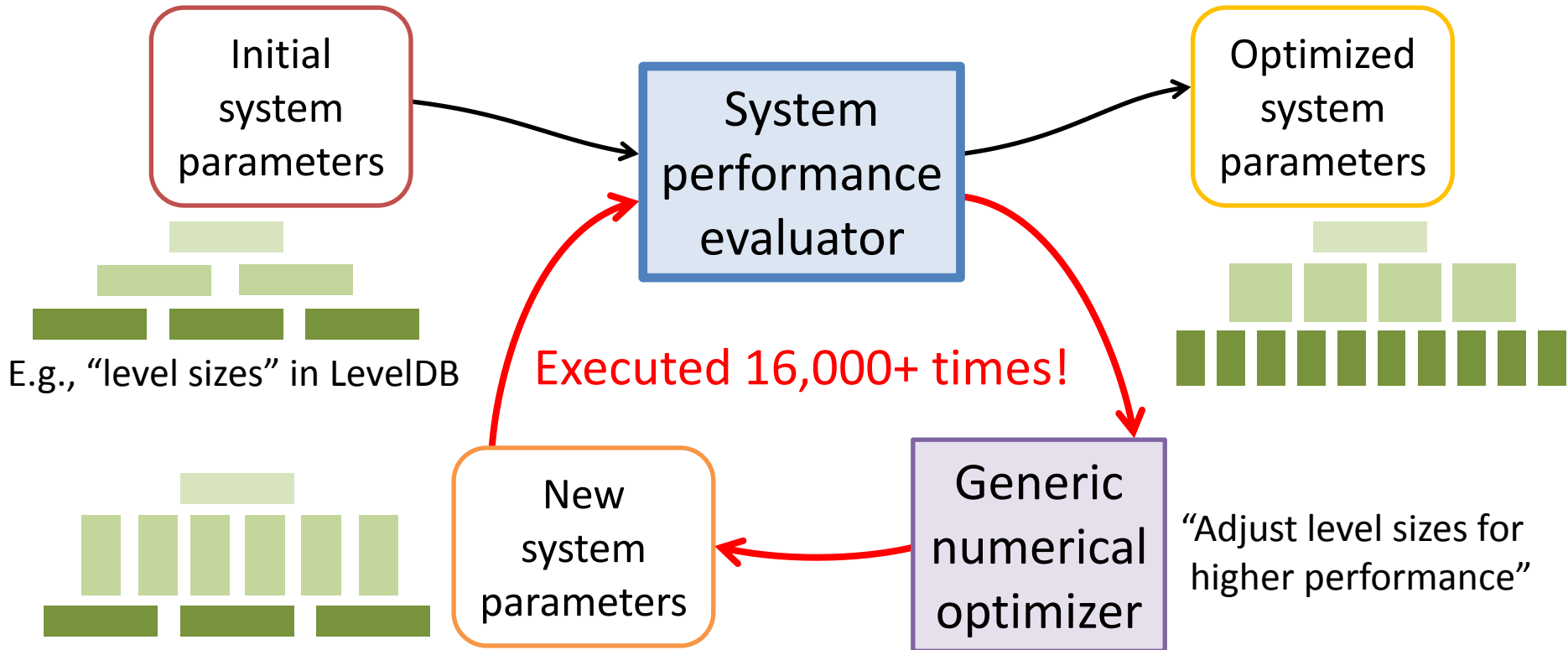
◇ **Asymptotic Analysis**
of core algorithms
(e.g., $O(\log N)$ I/Os per insert)

Want: **Accurate** and **fast** evaluation method

Experiment
using full implementation ◇
(e.g., 12 k inserts per second)

Accuracy

What You Can Do With Accurate and Fast Evaluation



Our level size optimization on LevelDB

- Up to 26.2% lower per-insert cost, w/o sacrificing query performance
- Finishes in 2 minutes (full experiment would take years)

Accurate and Fast Evaluation of MSLS Designs

Analytically model multi-stage log-structured designs using **new analytic primitives** that consider redundancy

Accuracy: Only $\leq 3\text{--}6.5\%$ off from LevelDB/RocksDB experiment

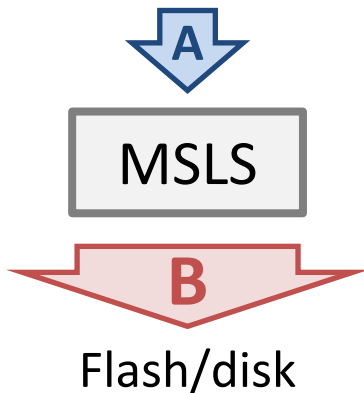
Speed: < 5 ms per run for a workload with 100 M unique keys

Performance Metric to Use

Focus of this talk: **Insert performance** of MSLS designs

- Often bottlenecked by writes to flash/disk
- Need to model amortized write I/O of inserts

User application

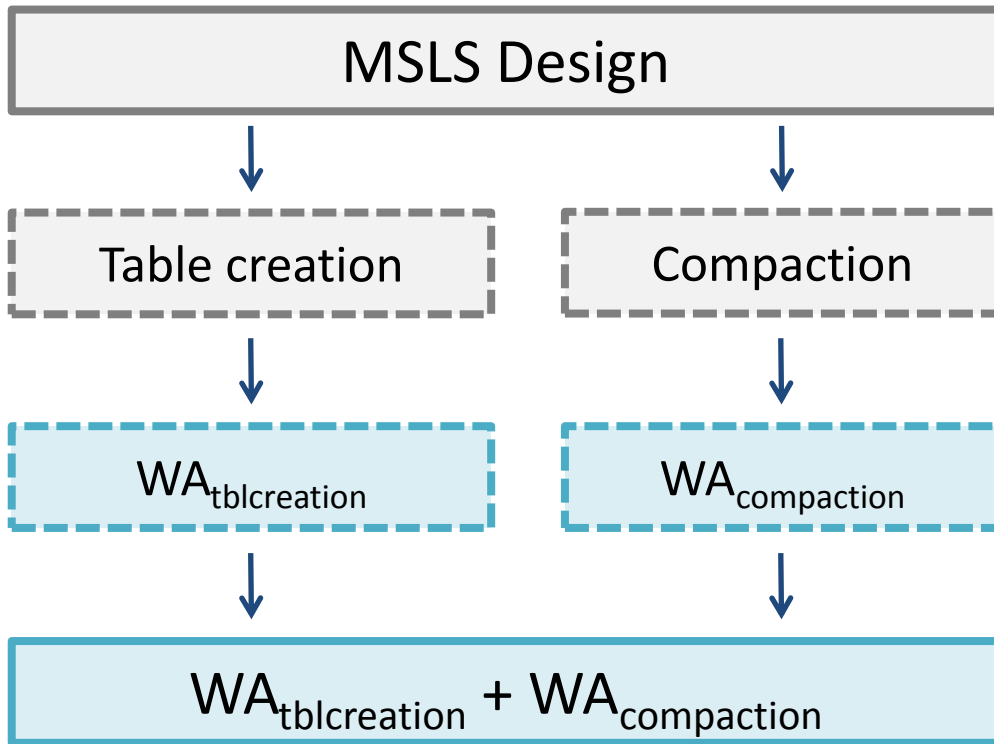


(Application-level) **Write amplification**

$$= \frac{\text{Size of data written to flash/disk (B)}}{\text{Size of inserted data (A)}}$$

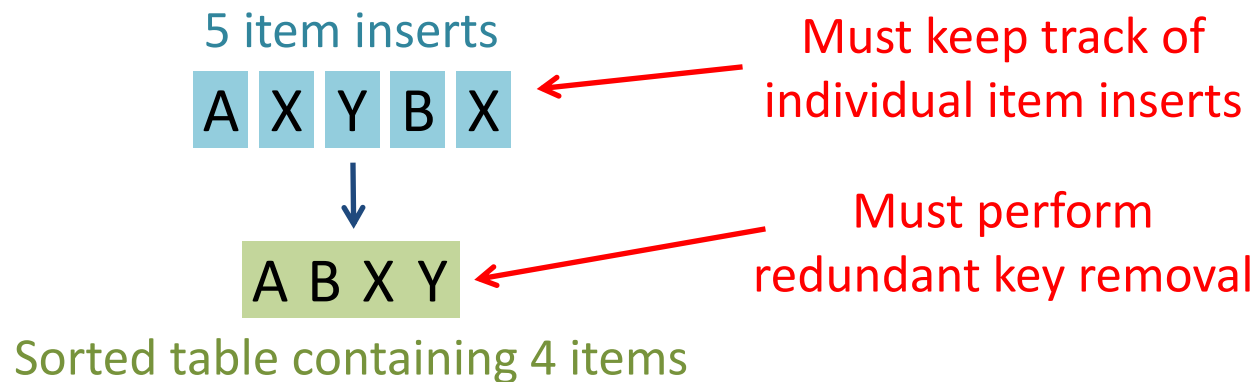
- Easier to analyze than raw throughput
- Closely related to raw throughput:
write amplification \propto 1/throughput

Divide-and-Conquer to Model MSLS Design



1. Break down MSLS design into small components
2. Model individual components' write amplification
3. Add all components' write amplification

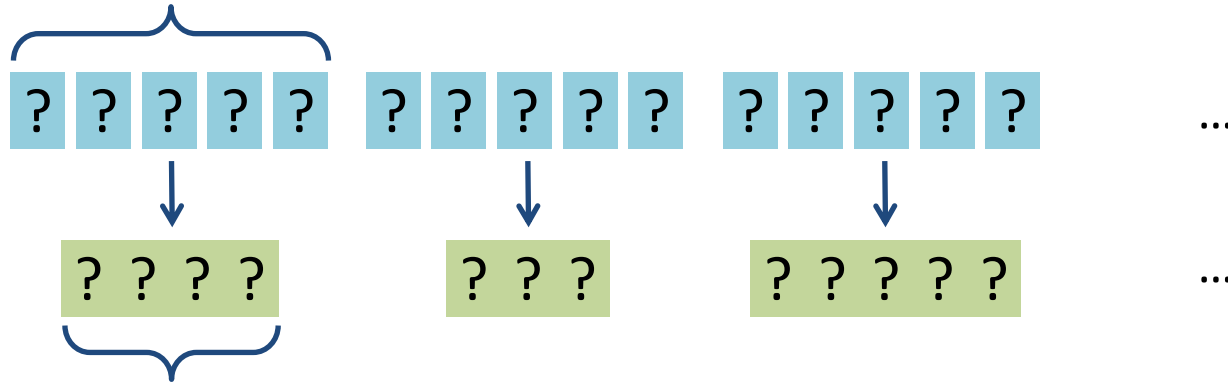
Modeling Cost of Table Creation: Strawman



Write amplification of this table creation event = $\frac{4}{5}$

Modeling Cost of Table Creation: Better Way

bufsize (max # of inserts buffered in memory)

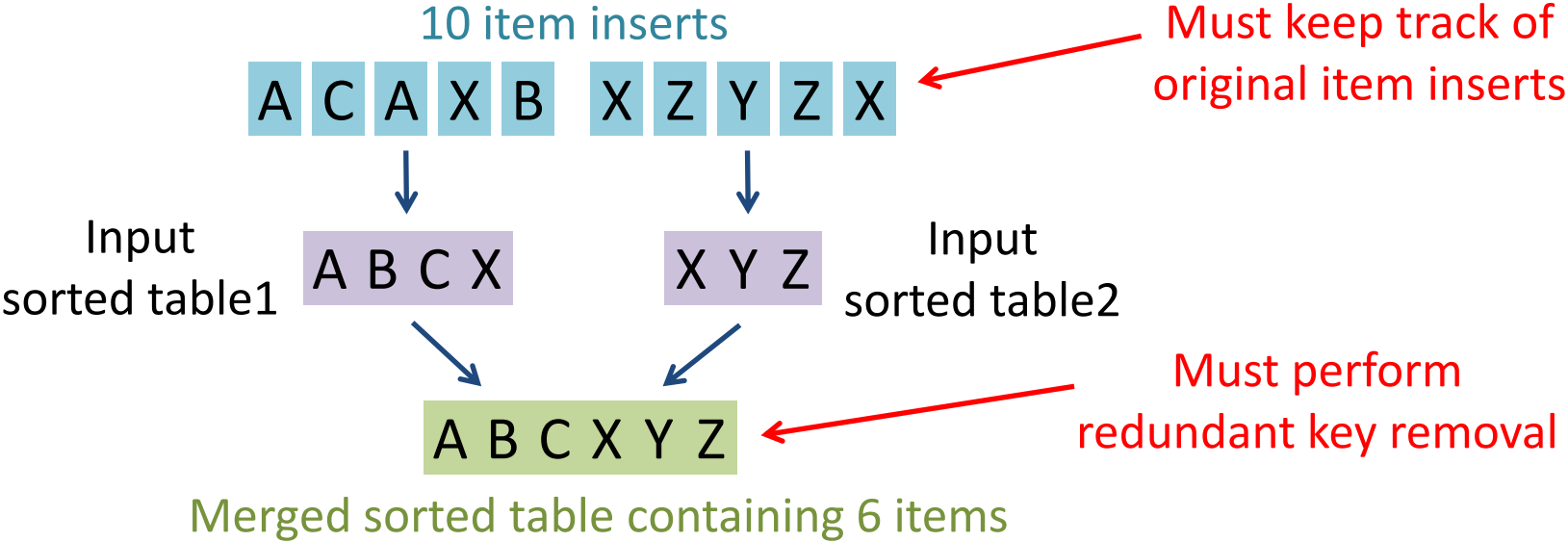


Unique(bufsize): expected # of **unique** keys in bufsize requests

$$\text{Write amplification of regular table creation} = \frac{\text{Unique}(\text{bufsize})}{\text{bufsize}}$$

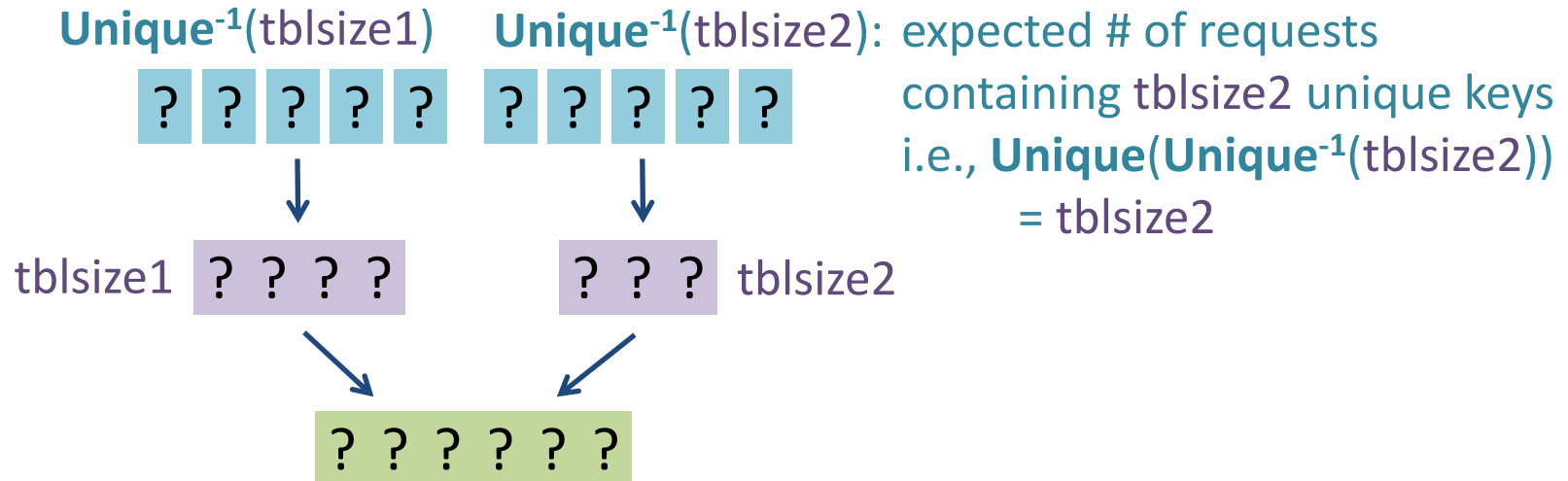
- ✓ No item-level information required
- ✓ Estimates general operation cost

Modeling Cost of Compaction: Strawman



Write amplification of this compaction event = $\frac{6}{10}$

Modeling Cost of Compaction: Better Way



$\text{Merge}(\text{tblsize1}, \text{tblsize2})$: expected # of unique keys in input tables whose sizes are tblsize1 and tblsize2

Write amplification of 2-way compaction =
$$\frac{\text{Merge}(\text{tblsize1}, \text{tblsize2})}{\text{Unique}^{-1}(\text{tblsize1}) + \text{Unique}^{-1}(\text{tblsize2})}$$

- ✓ No item-level information required
- ✓ Estimates general operation cost

New Analytic Primitives Capturing Redundancy

Unique: [# of requests] \rightarrow [# of unique keys]

Unique⁻¹: [# of requests] \leftarrow [# of unique keys]

Merge: [multiple # of unique keys] \rightarrow [total # of unique keys]

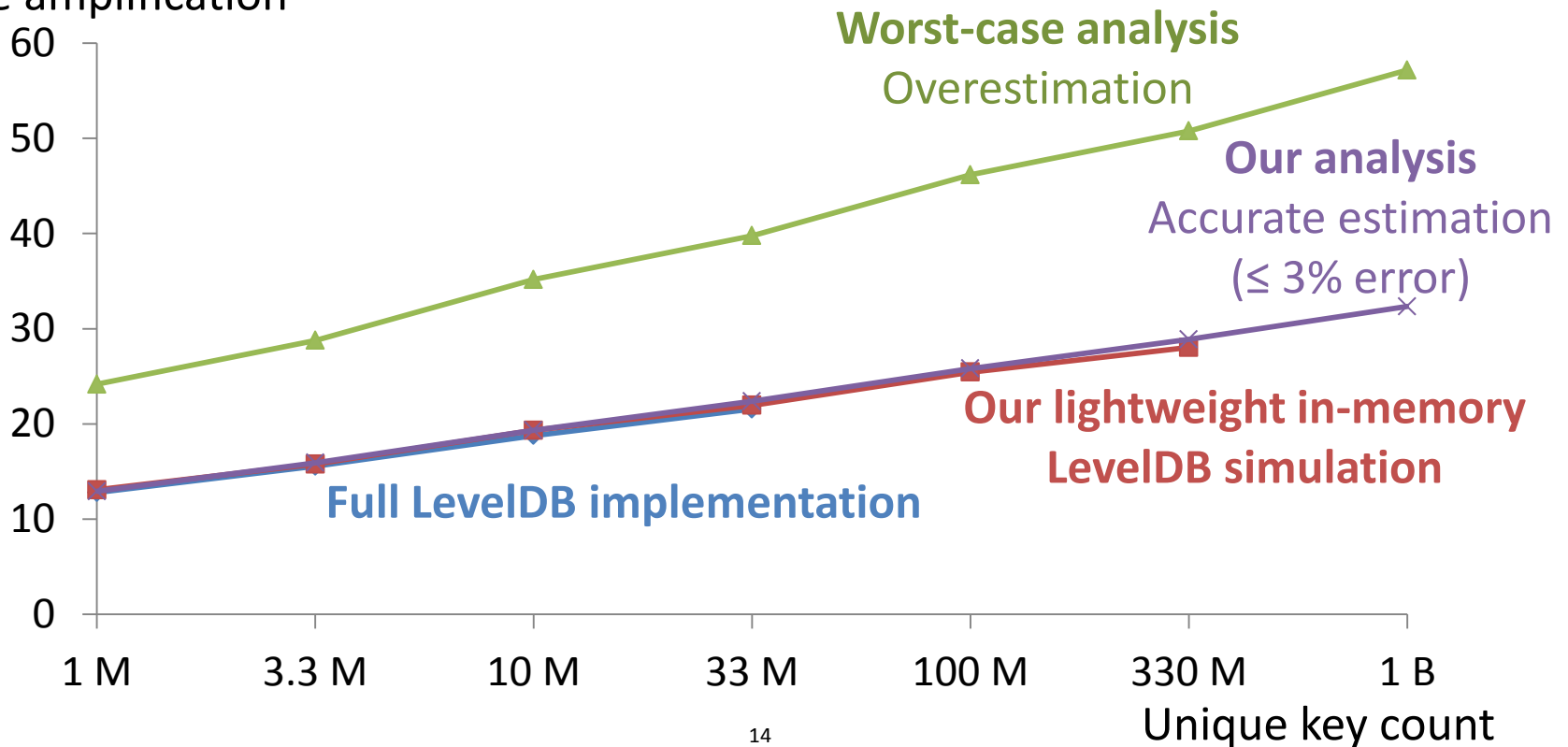
- **Fast** to compute (see paper for mathematical descriptions)
- Consider **redundancy**: $\text{Unique}(p) \leq p$ $\text{Merge}(u, v) \leq u + v$
- Reflect **workload skew**: $[\text{Unique}(p) \text{ for Zipf}] \leq [\text{Unique}(p) \text{ for uniform}]$
- **Caveat**: Assume no or little dependence between requests

High Accuracy of Our Evaluation Method

Compare measured/estimated write amplification of insert requests on LevelDB

- Key-value item size: 1,000 bytes
- Unique key count: 1 million–1 billion (1 GB–1 TB)
- Key popularity dist.: Uniform

Write amplification



High Speed of Our Evaluation Method

Compare **single-run** time to obtain write amplification of insert requests for a **specific** workload using a **single** set of system parameters

- LevelDB implementation: fsync disabled
- LevelDB simulation: in-memory, optimized for insert processing

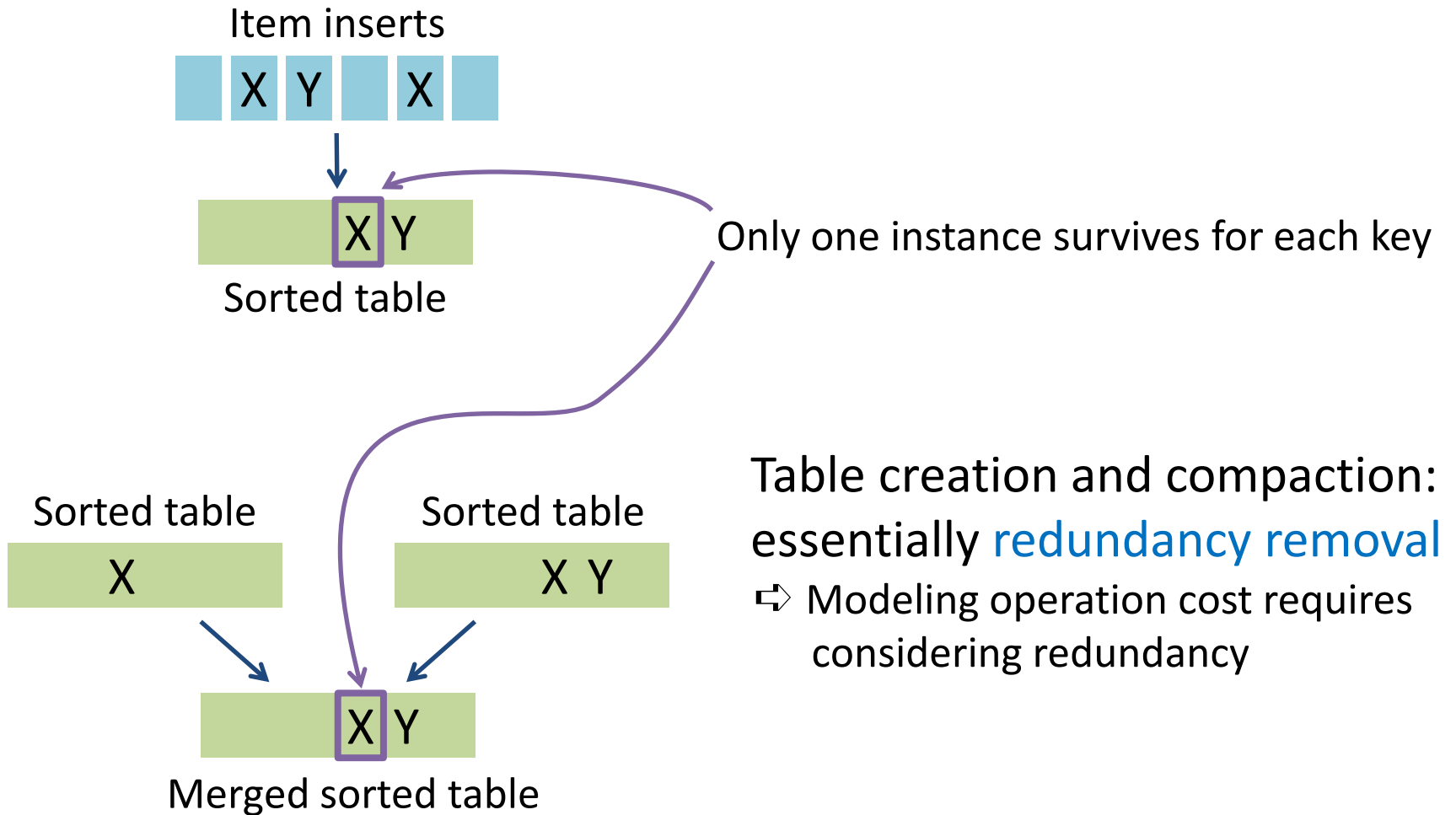
Method	Workload size (# of unique keys)	Elapsed time
Experiment using LevelDB implementation	10 M	101 minutes
Experiment using LevelDB simulation	100 M	45 minutes
Our analysis	100 M	< 5 ms

Summary

- Evaluation method for multi-stage log-structured designs
 - **New analytic primitives** that consider redundancy
 - **System models** using new analytic primitives
- **Accurate and fast**
 - Only $\leq 3\text{--}6.5\%$ error in estimating insert cost of LevelDB/RocksDB
 - **Several orders of magnitude faster** than experiment
- Example applications
 - **Automatic system optimization** ($\sim 26.2\%$ faster inserts on LevelDB)
 - **Design improvement** ($\sim 32.0\%$ faster inserts on RocksDB)
- Code: github.com/efficient/msls-eval

Backup Slides

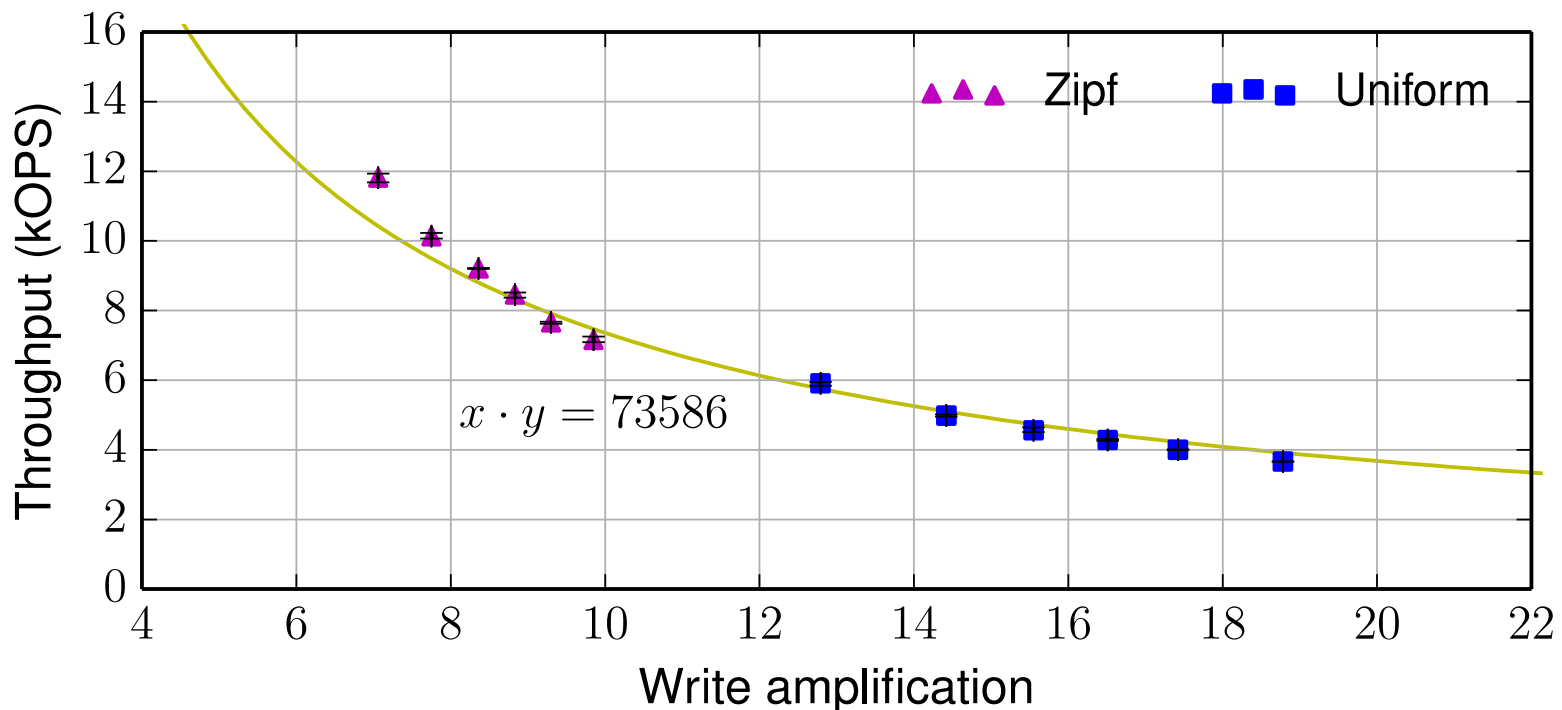
Nature of MSLS Operations



Write Amplification vs. Throughput

Compare measured write amplification/throughput of insert requests on LevelDB

- Key-value item size: 1,000 bytes
- Unique key count: 1 million–10 million (1 GB–10 GB)
- Key popularity dist.: Uniform, Zipf (skew=0.99)



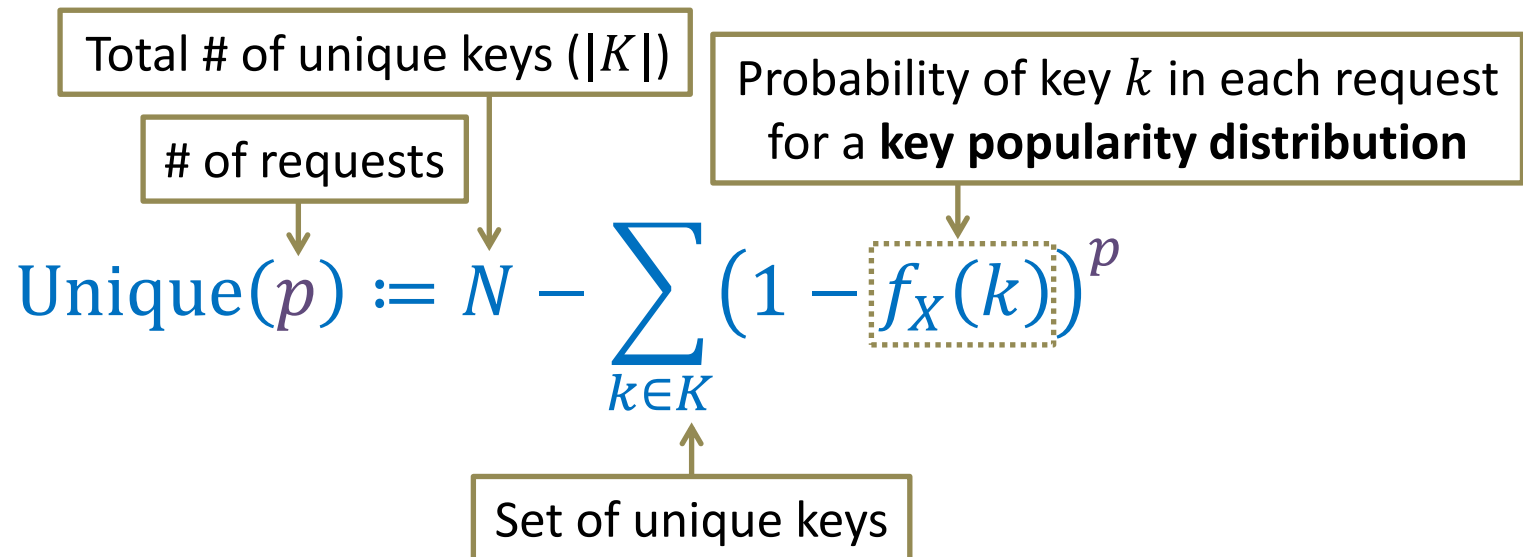
Mathematical Description of New Primitives

Unique: [# of requests] \rightarrow [# of unique keys]

Unique⁻¹: [# of requests] \leftarrow [# of unique keys]

Merge: [multiple # of unique keys] \rightarrow [total # of unique keys]

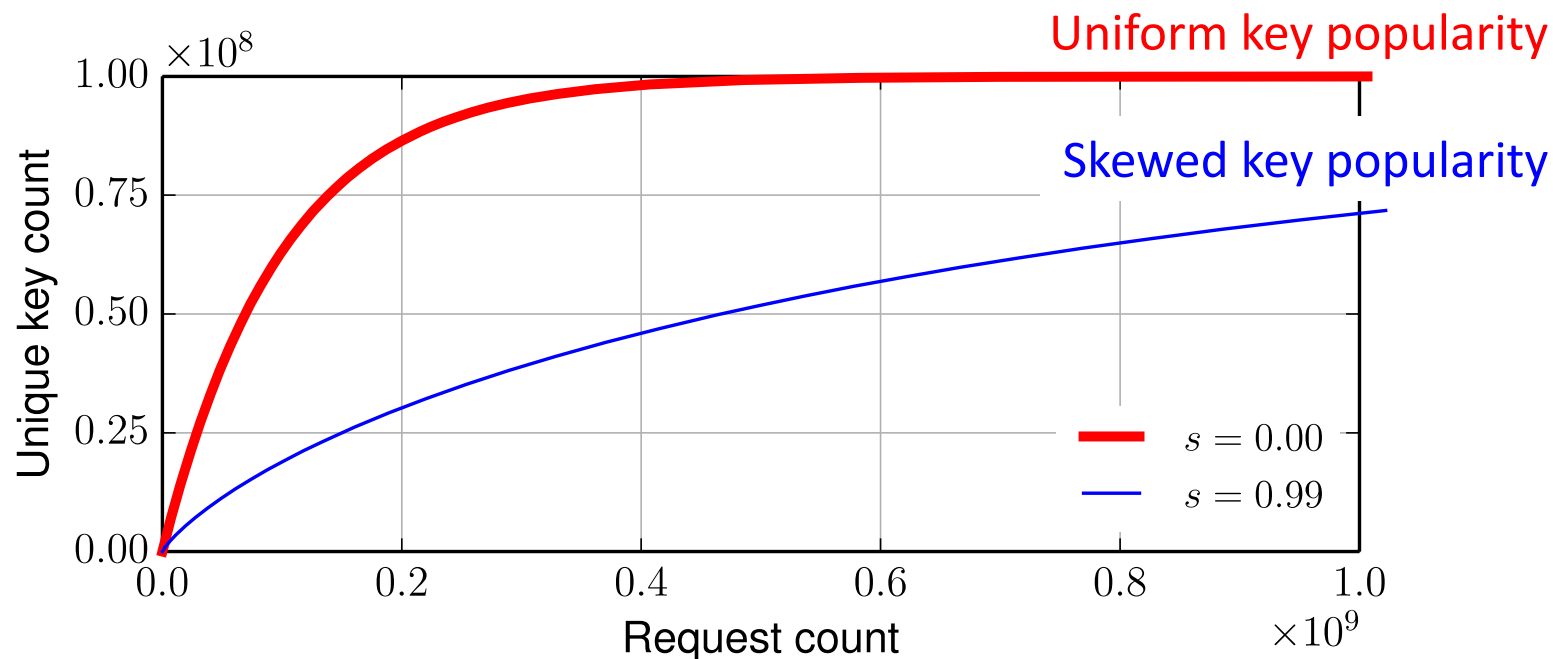
$$\text{Merge}(u, v) = \text{Unique}(\text{Unique}^{-1}(u) + \text{Unique}^{-1}(v))$$



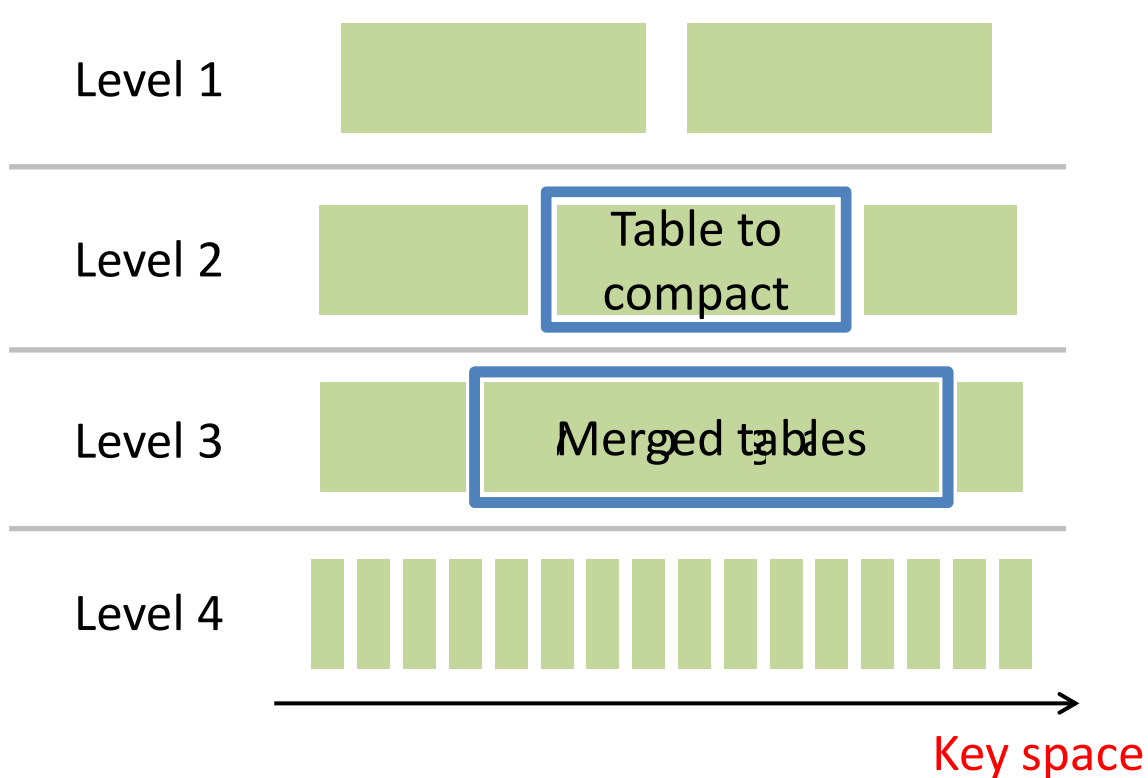
Unique as a Function of Request Count

Compare measured write amplification/throughput of insert requests on LevelDB

- Key-value item size: 1,000 bytes
- Unique key count: 100 M (100 GB)
- Request count: 0–1 billion
- Key popularity dist.: Uniform, Zipf (skew=0.99)



LevelDB Design Overview



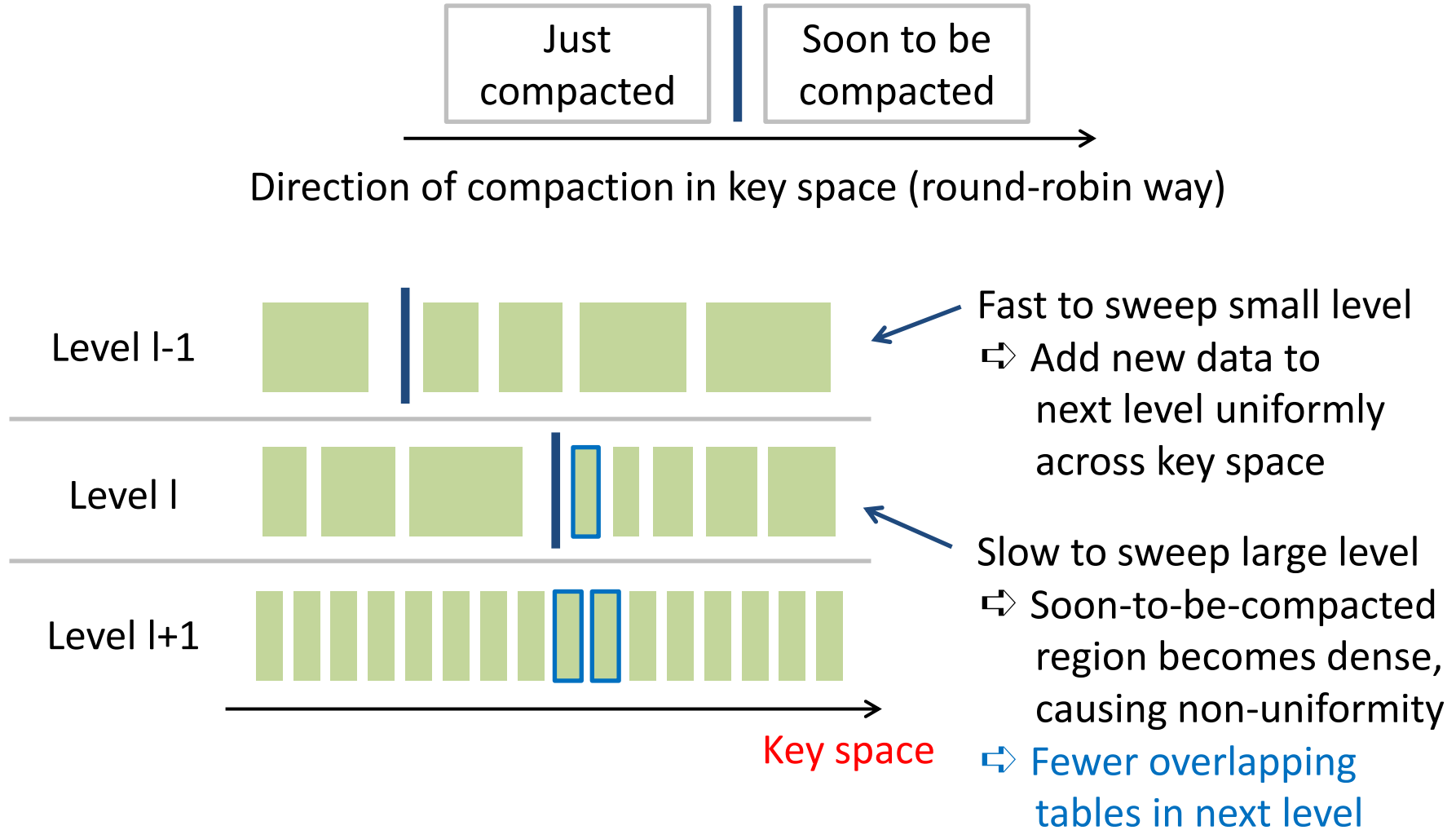
Each level's total size
= ~**10X** previous level's

Each level are partitioned
into small tables (~2 MB)
for **incremental compaction**

Q: Average # of overlaps?
⇒ Less than **10!**
("non-uniformity")

(Omitted: memtable, write-ahead log, level 0)

Non-Uniformity in LevelDB

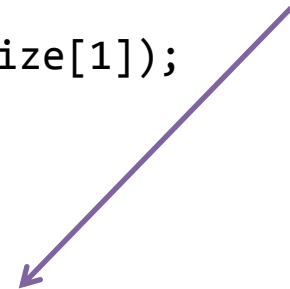


(Omitted: memtable, write-ahead log, level 0)

Pseudo Code of LevelDB Model

```
1 // @param L      maximum level
2 // @param wal    write-ahead log file size
3 // @param c0     level-0 SSTable count
4 // @param size   level sizes
5 // @return      write amplification (per-insert cost)
6 function estimateWA_LevelDB(L, wal, c0, size[]) {
7     local l, WA, interval[], write[];
8
9     // mem -> log
10    WA = 1;
11
12    // mem -> level-0
13    WA += unique(wal) / wal;
14
15    // level-0 -> level-1
16    interval[0] = wal * c0;
17    write[1] = merge(unique(interval[0]), size[1]);
18    WA += write[1] / interval[0];
19
20    // level-1 -> level-(l+1)
21    for (l = 1; l < L; l++) {
22        interval[l] = interval[l-1] + dinterval(size, l);
23        write[l+1] = merge(unique(interval[l]), size[l+1]) + unique(interval[l]);
24        WA += write[l+1] / interval[l];
25    }
26
27    return WA;
28 }
```

LevelDB-specific function
to take into account “non-uniformity”

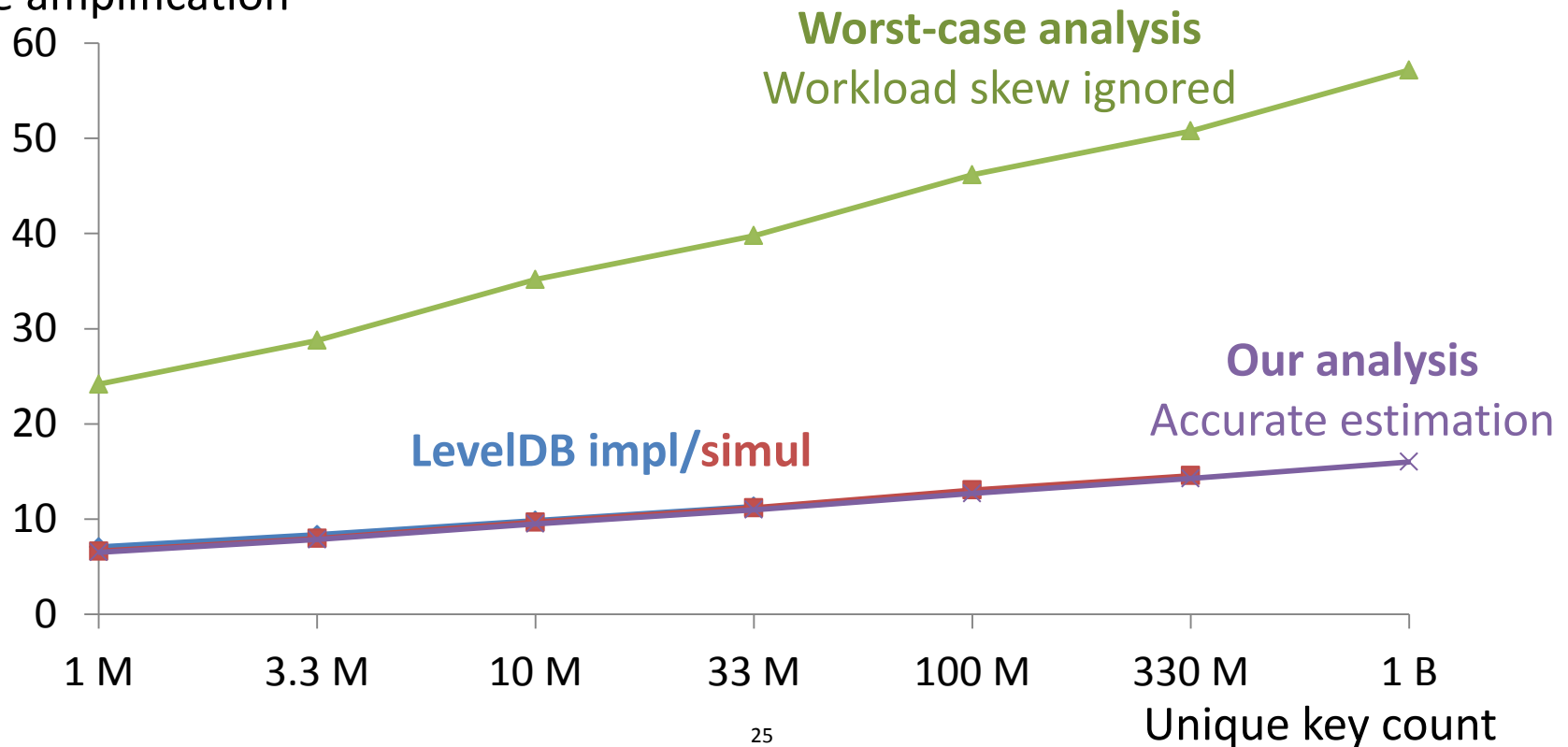


Sensitivity to Workload Skew

Compare measured/estimated write amplification of insert requests on LevelDB

- Key-value item size: 1,000 bytes
- Unique key count: 1 million–1 billion (1 GB–1 TB)
- **Key popularity dist.: Zipf (skew=0.99)**

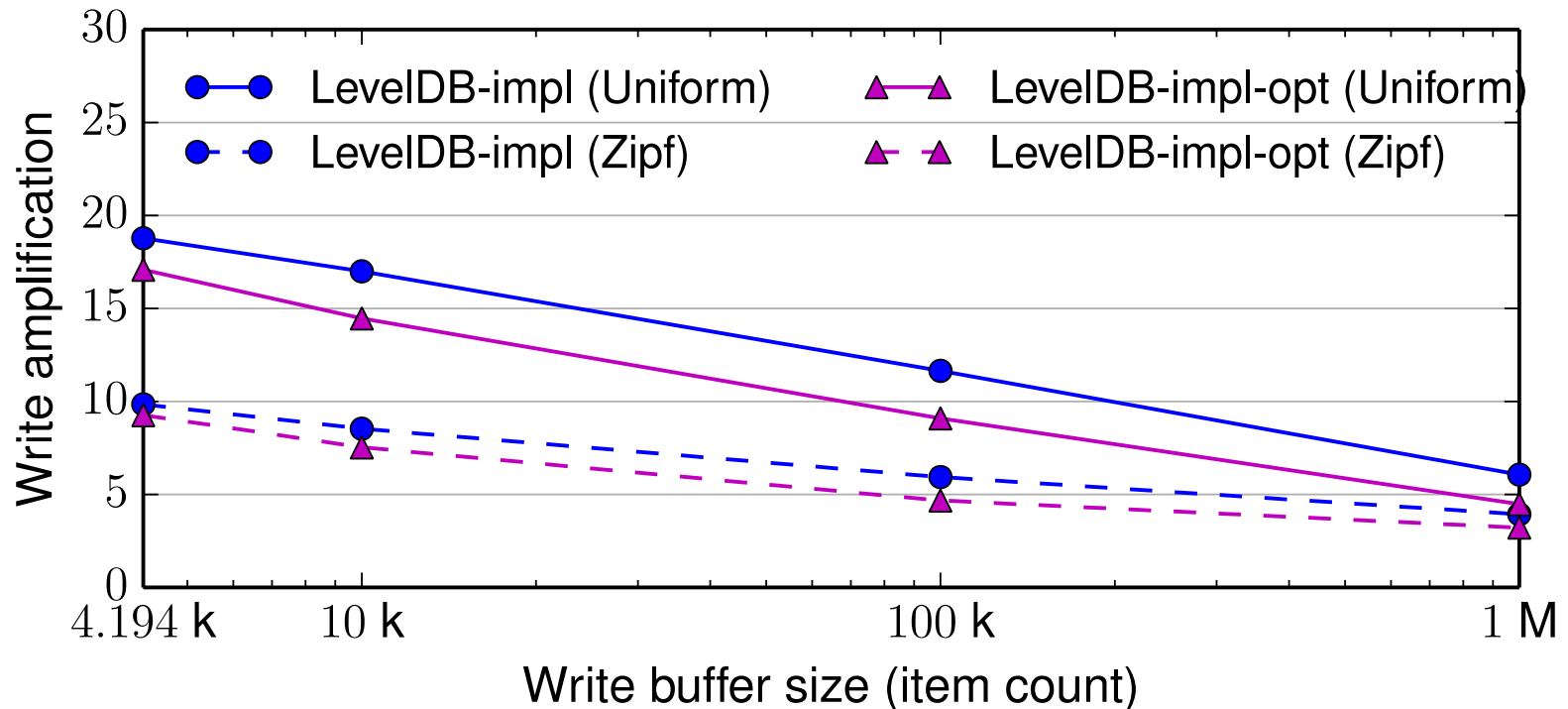
Write amplification



Automatic System Optimization Result

Compare measured/estimated write amplification of insert requests on LevelDB

- Key-value item size: 1,000 bytes
- Write buffer size: 4 MiB–[10% of total unique key count]
- Unique key count: 10 million (10 GB)
- Key popularity dist.: Uniform, Zipf (skew=0.99)



End of Slides
