# Thomas Kissinger

*DEXTER*
*Parallel In-Memory Indexing and*
*Direct Query Processing on Prefix Trees*

> **Application-Level Trends**

Database **Technology**
Group

- *Evolution of DWH Applications*

| **Reporting** | **Analysis** | **Advanced Analytics** | **Operational BI** |
|---|---|---|---|
| *What did happen?* | *Why did it happen?* | *What will happen?* | *What happens right now?* |
| *Create reports with pre-defined queries.* | *Increasing number of ad-hoc queries.* | *Extension of the analytical model (e.g. Forecasting).* | *Continuous streams of ad-hoc queries and propagated updates.* |
| Step 1 | Step 2 | Step 3 | Step 4 |

- Batch
- Adhoc
- Analytics
- Updates

- *Advanced Analytics*
  - Sophisticated statistical models
  - Machine learning
  - → Mixed worklods

- *Operational BI*
  - →High update rates
  - →Transactional workloads

- *SAP HANA, C-Store, ...*

ad-hoc queries

periodic merge

write-optimized
temporary store

update flow

read-optimized column stores
- In-memory
    - About 30GB/s theoretical peak read performance
- Highly parallel
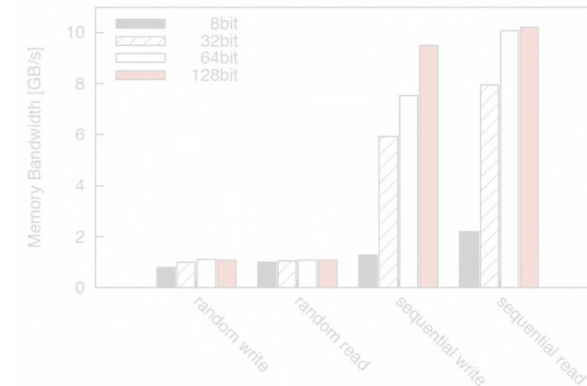
- *Shortcomings*
    - Poor point query support (e.g. Advanced Analytics)
    - Additional query of temporary store necessary
    - Periodic merge rebuilds column store and stalls queries in the meantime

- *Shifting Memory Hierarchy*
  - In-memory databases/indexing
  - Cache-awareness
  - Other Block/MTU sizes
  - Sequential access still faster than random access



- *More and more Cores*
  - clock rate got stuck
  - Massive parallelism
  - Optimized software required

I7-2600
8 HW Threads

Xeon E7
20 HW Threads

Knight's Ferry
128 HW Threads

- *Modern Hardware Architectures*
  - SIMD instructions
  - Multiple memory channels
  - SMP changes to NUMA
  - Heterogeneous Hardware (GPUs, Co-processors …)
  - Adaptive Hardware (FPGAs, …)

▪ **FAST: fast architecture sensitive tree search on modern CPUs and GPUs**

➢ **SIGMOD 2010**

*Changkyu Kim, Jatin Chhugani, Nadathur Satish, Eric Sedlar, Anthony D. Nguyen, Tim Kaldewey, Victor W. Lee, Scott A. Brandt and Pradeep Dubey*

• Memory hierarchy optimized binary tree
• Read-optimized index structure
• 51M reads/s, but only 10 updates/s for a tree of 64M keys

▪ **PALM: Parallel Architecture-Friendly Latch-Free Modifications to B+ Trees on Many-Core Processors**

➢ **VLDB 2011**

*Jason Sewall, Jatin Chhugani, Changkyu Kim, Nadathur Satish and Pradeep Dubey*

• B+-Tree based index
• Synchronous batch updates to avoid latches

Both index structures suffer from a poor update performance

❖ INTRODUCTION AND TRENDS

❖ DEXTER PROJECT
- ◆ Overview
- ◆ Project Map

❖ CORE INDEXING

❖ DIRECT QUERY PROCESSING

❖ CONCLUSIONS

*(**Dr**esden Inde**x** for **T**ransactional Access on **E**me**r**ging Technologies)*

▪ *Transactional Index, optimized for mixed queries and high update rates*

▪ *Sponsored by SFB 912: HAEC (Highly Adaptive Energy-Efficient Computing)*

Matthias Böhm

Thomas Kissinger

Peter B. Volk
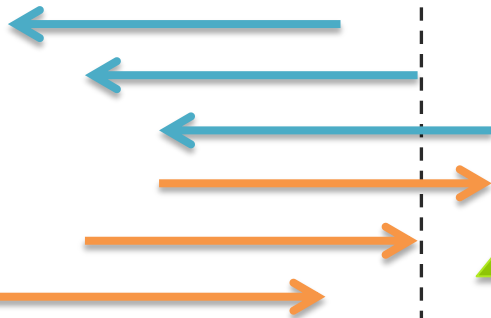
Benjamin Schlegel

Ulrike Fischer

Dirk Habich

DEXTER

*Programming Contest*

**2nd place 2009**
indexing system for main memory data

**1st place 2011**
high-throughput main-memory index which is made durable using a flash-based SSD

*Direct Query Processing*

high point query throughput

*Core Indexing*

high update rates

❖ INTRODUCTION AND TRENDS

❖ DEXTER PROJECT

❖ CORE INDEXING

    ◆ Generalized Prefix Trees
    ◆ Design Space for parallel In-Memory Indexing

❖ DIRECT QUERY PROCESSING

❖ CONCLUSIONS

| 0000 | 0000 | 0110 | 1011 |
|------|------|------|------|
| 0 | 0 | 6 | 11 |

Static Prefix Length: **k' = 4**

- Only one 64bit access per node
- Deterministic path
- No Balancing

Static Prefix Length: **k' = 1**

Static Prefix Length: **k' = 4**

Static Prefix Length: **k' = 16**



flatter tree → Less memory transfers

deeper tree → More memory transfers / Better memory utilization

**Variable prefix length on tree level or node level possible**

Database Technology
Group

- *Two important dimensions for parallel in-memory indexing*

|  | No | Batching | Yes |
|---|---|---|---|
|  | | | |

No — FAST — DEXTER

No — DEXTER Batched Reads — ➕ Direct Processing — ➖ Synchronization

Partitioning

Yes — (hatched) — PALM — DEXTER BUZZARD — ➕ Independent Data — ➖ Indirect Processing

➕ Low latency  ➕ High Throughput

➖ Increased Latency

▪ *Two important dimensions for parallel in-memory indexing*

|  | No | Batching | Yes |
|---|---|---|---|
| **No** | FAST / DEXTER | | DEXTER Batched Reads |
| **Partitioning** | | | |
| **Yes** | | | PALM / DEXTER BUZZARD |

➕ Direct Processing

➖ Synchronization

➕ Independent Data

➖ Indirect Processing

➕ Low latency  ➕ High Throughput

➖ Increased Latency

Read-Copy-Updates (RCU) - aware Memory Manager

- Heavyweight Latches
  - Mutex/Futex
  - Spinlock

- Lightweight Atomics
  - Compare-and-Swap (CAS)

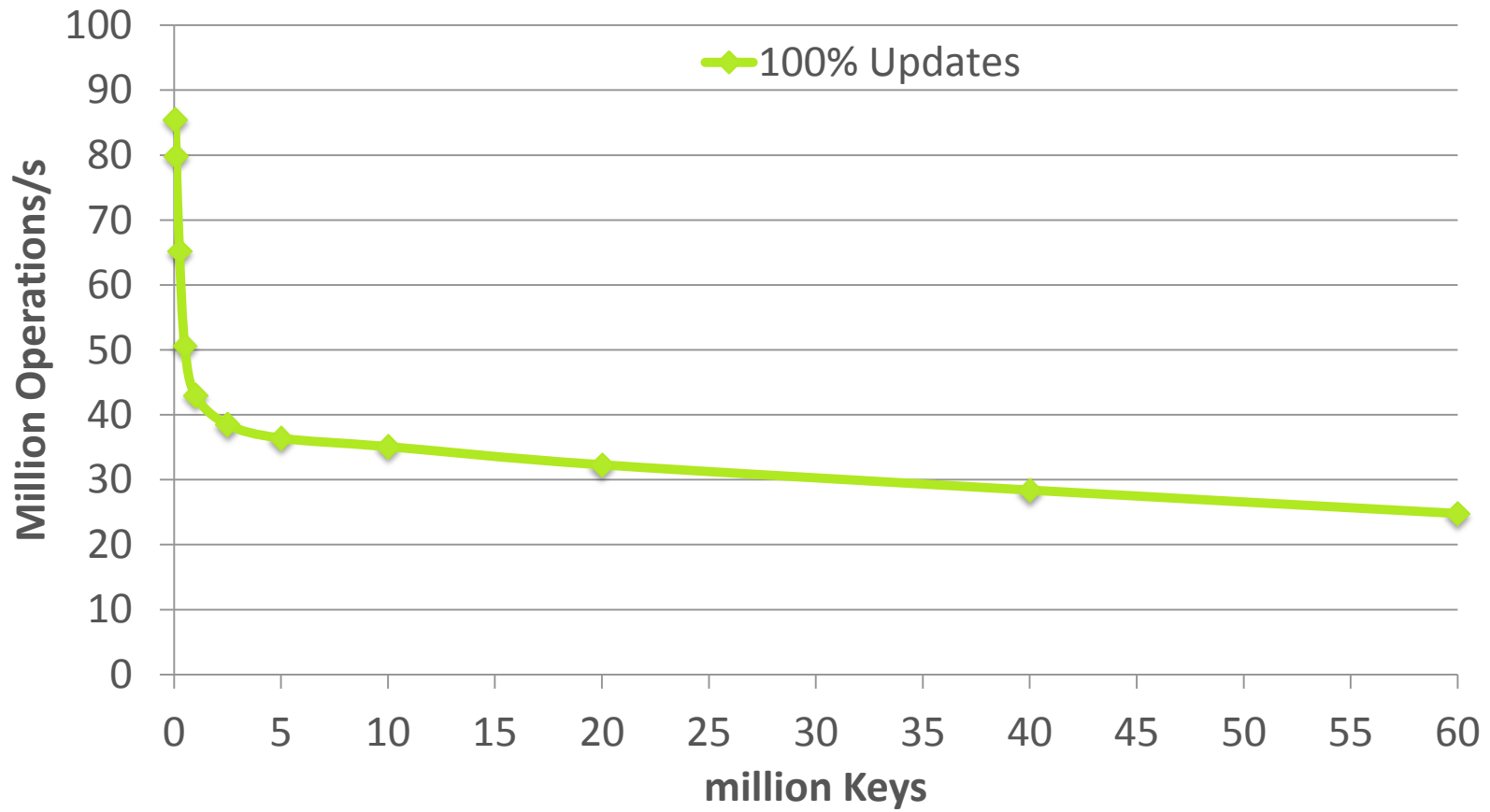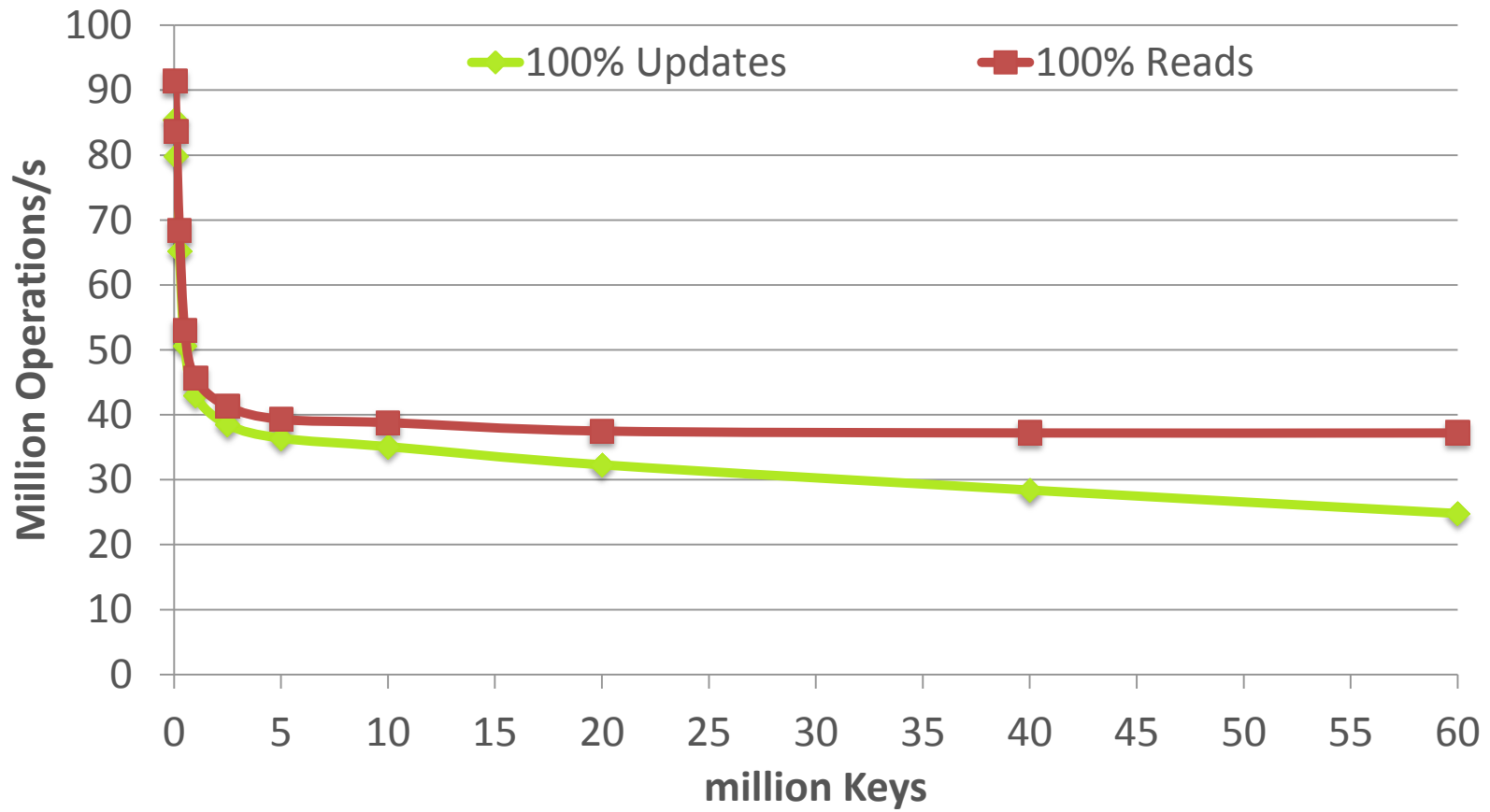| 0000 0 | 0000 0 | 0110 6 | 1011 11 |



**Read-Copy-Updates (RCU) - aware Memory Manager**

■ *Intel i7-2600 (4 cores @3.4GHz, 2 memory channels, 16GB DDR3 @1333MHz)*

- *Intel i7-2600 (4 cores @3.4GHz, 2 memory channels, 16GB DDR3 @1333MHz)*
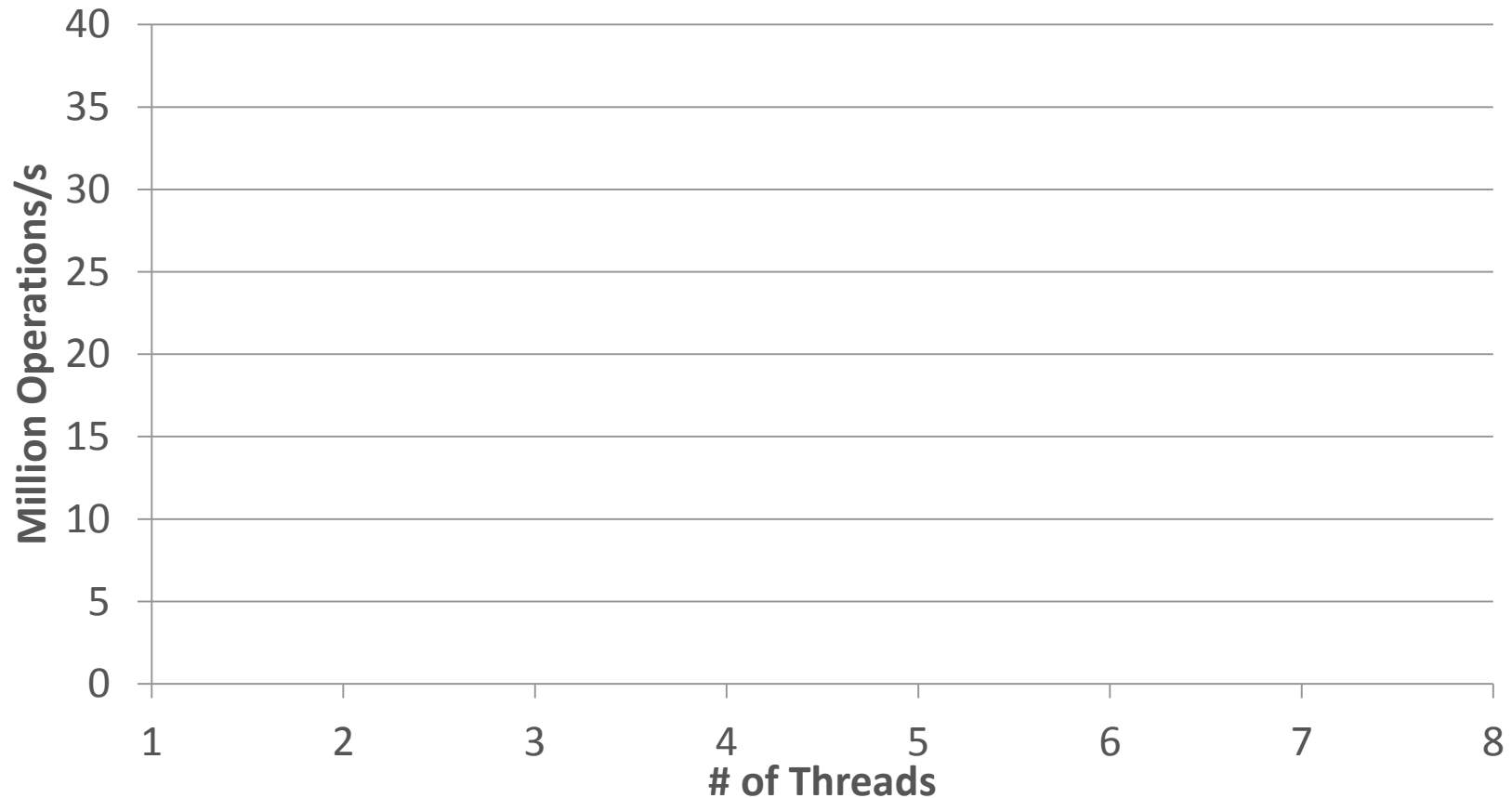
Database Technology
Group

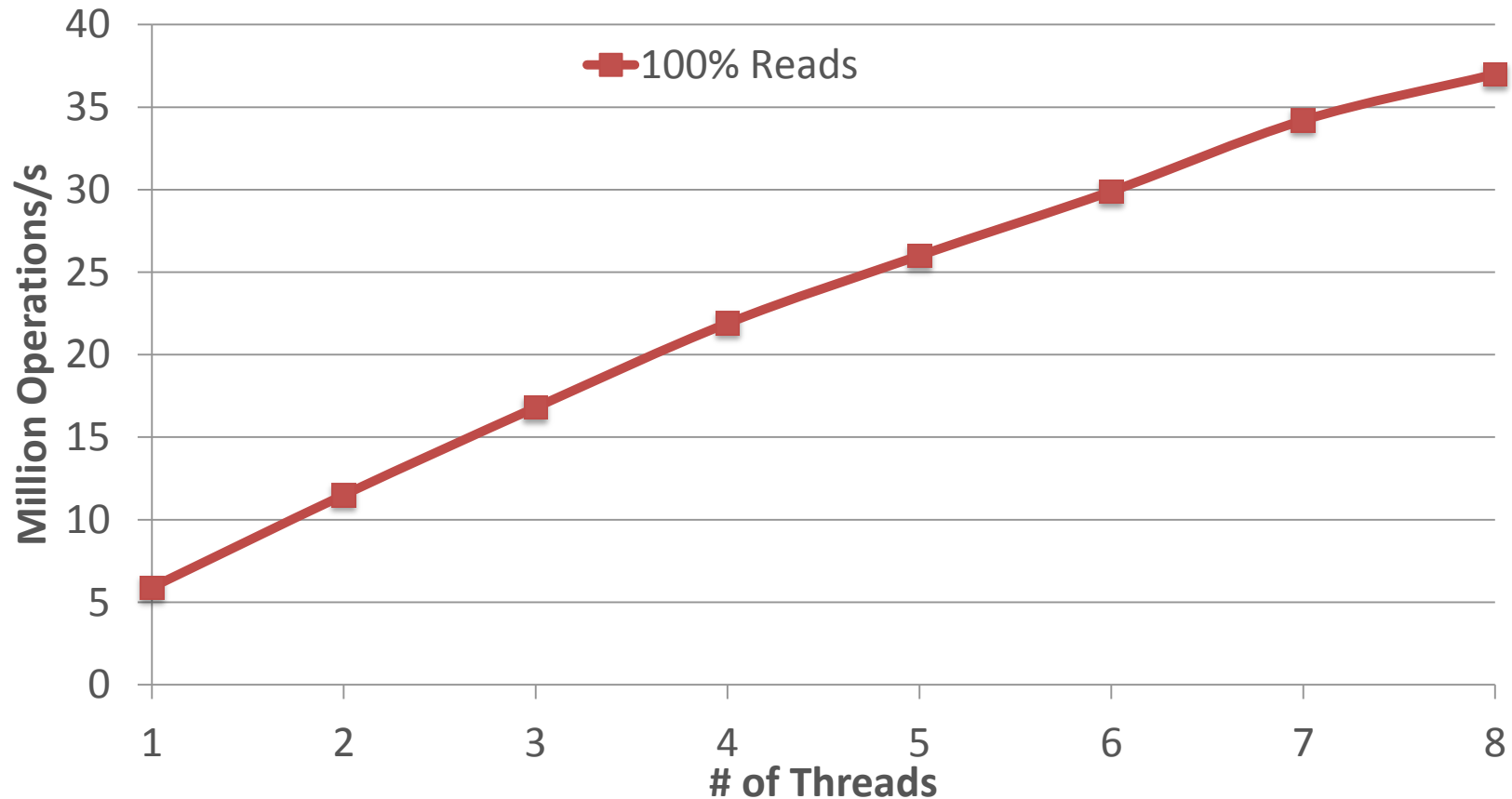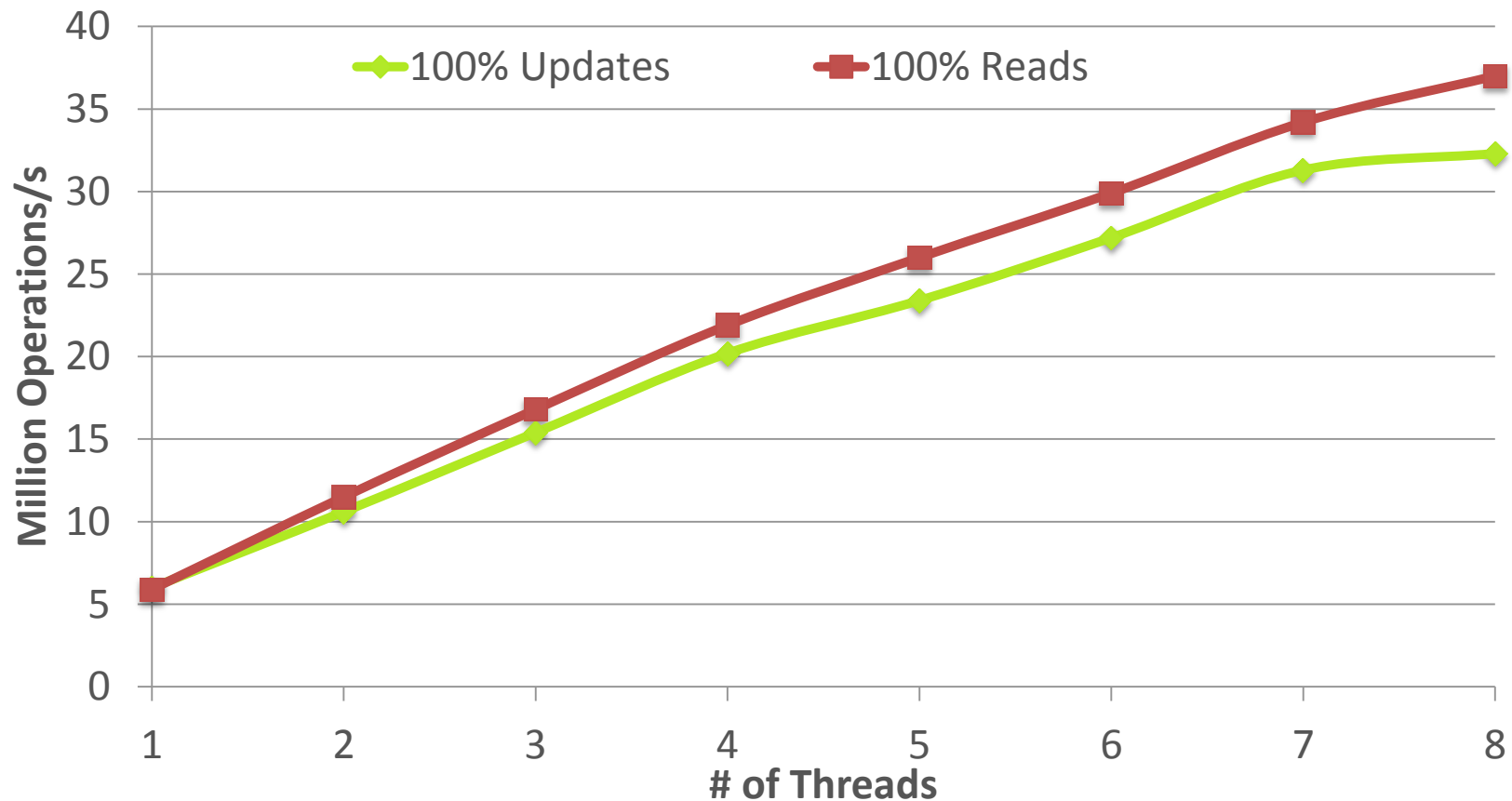- *Intel i7-2600 (4 cores @3.4GHz, 2 memory channels, 16GB DDR3 @1333MHz)*

- *Intel i7-2600 (4 cores @3.4GHz, 2 memory channels, 16GB DDR3 @1333MHz)*
- *20M Keys (32bit)*

- *Intel i7-2600 (4 cores @3.4GHz, 2 memory channels, 16GB DDR3 @1333MHz)*
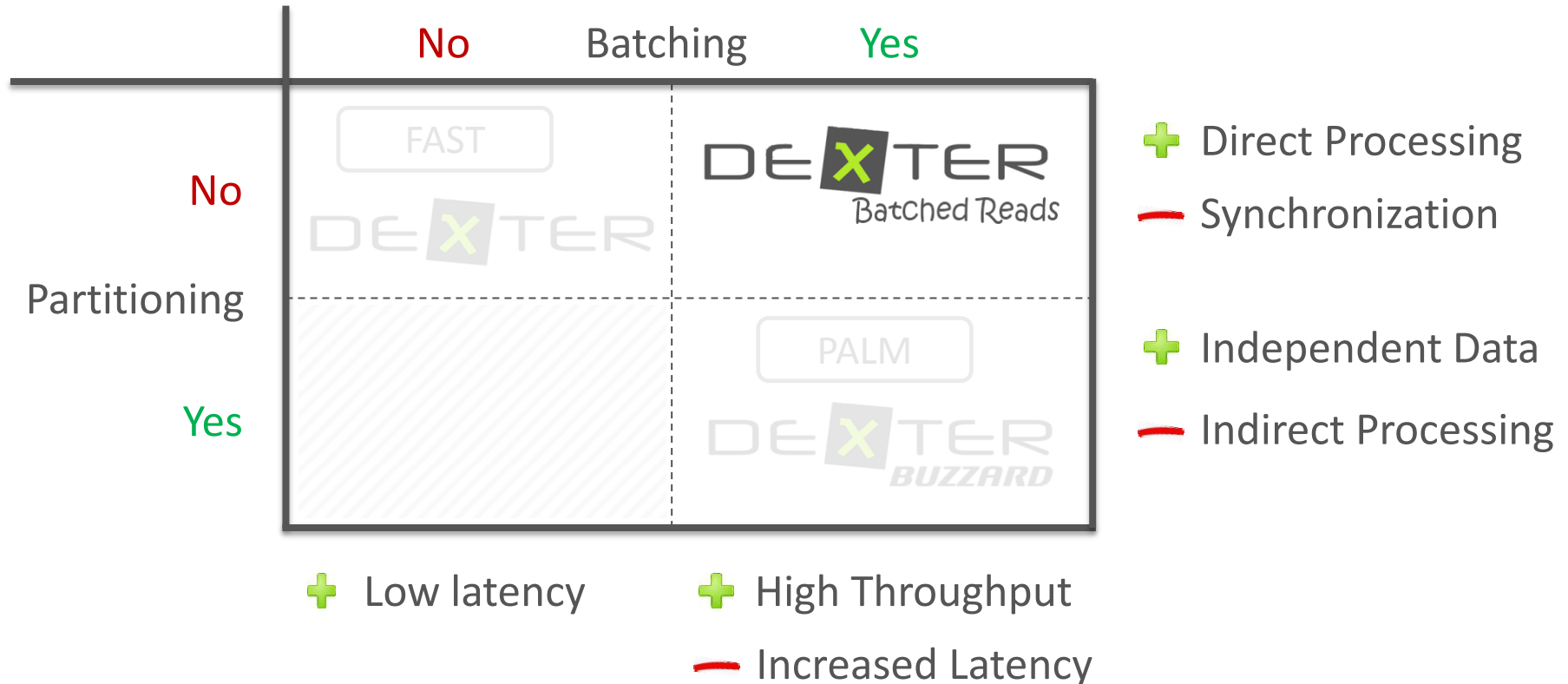- *20M Keys (32bit)*

- *Intel i7-2600 (4 cores @3.4GHz, 2 memory channels, 16GB DDR3 @1333MHz)*
- *20M Keys (32bit)*

■ *Two important dimensions for parallel in-memory indexing*

|  | **No** Batching **Yes** |  |
|---|---|---|
| **No** | FAST DEXTER | DEXTER *Batched Reads* |
| Partitioning **Yes** |  | PALM DEXTER *BUZZARD* |

➕ Direct Processing

➖ Synchronization

➕ Independent Data

➖ Indirect Processing

➕ Low latency    ➕ High Throughput

➖ Increased Latency

Key Batch:

| 107 | 111 | 8874 | 60928 |

Batch Size = 4



- Reads some nodes (especially top-level nodes) only once
- Prefetching possible
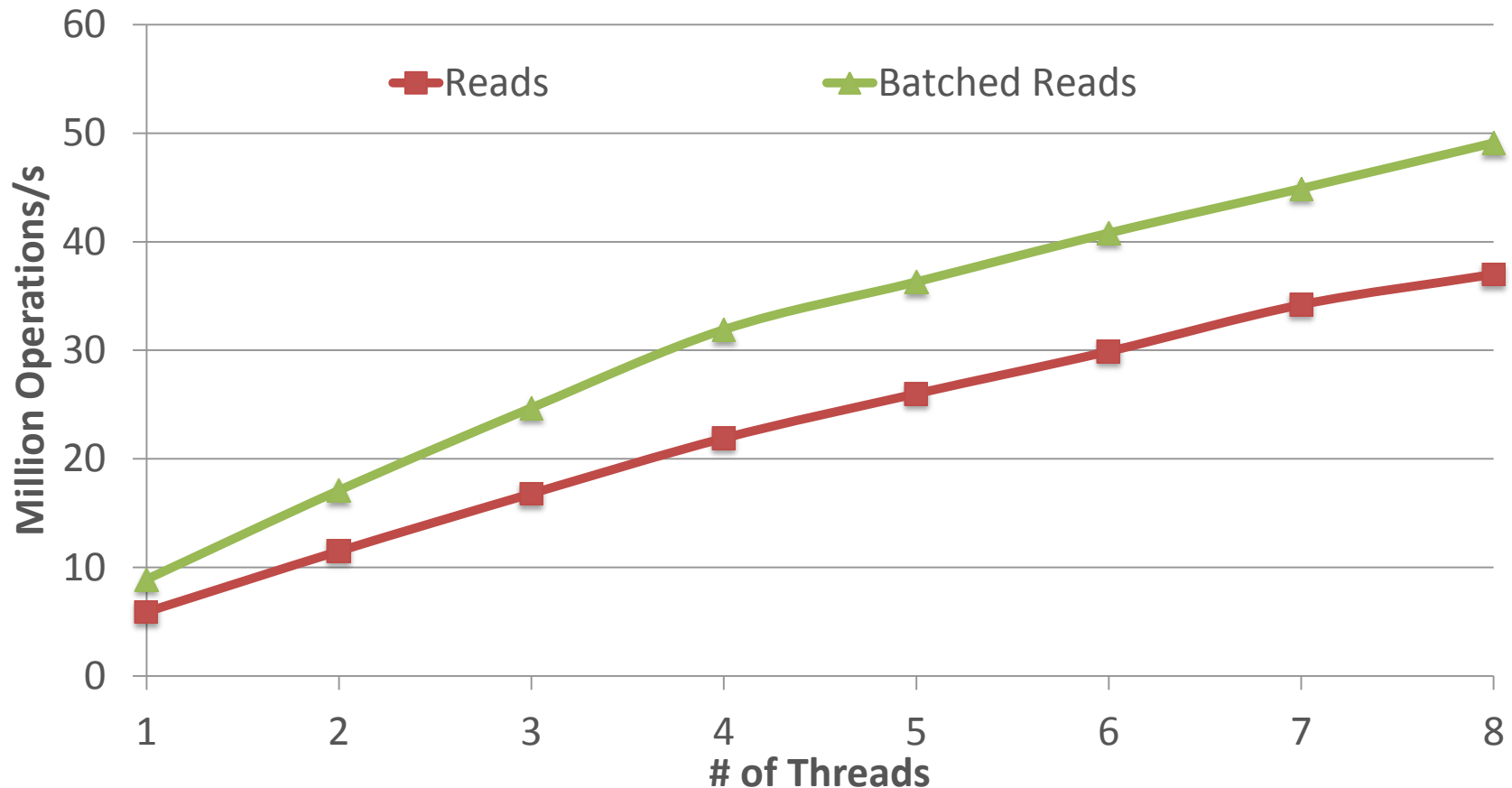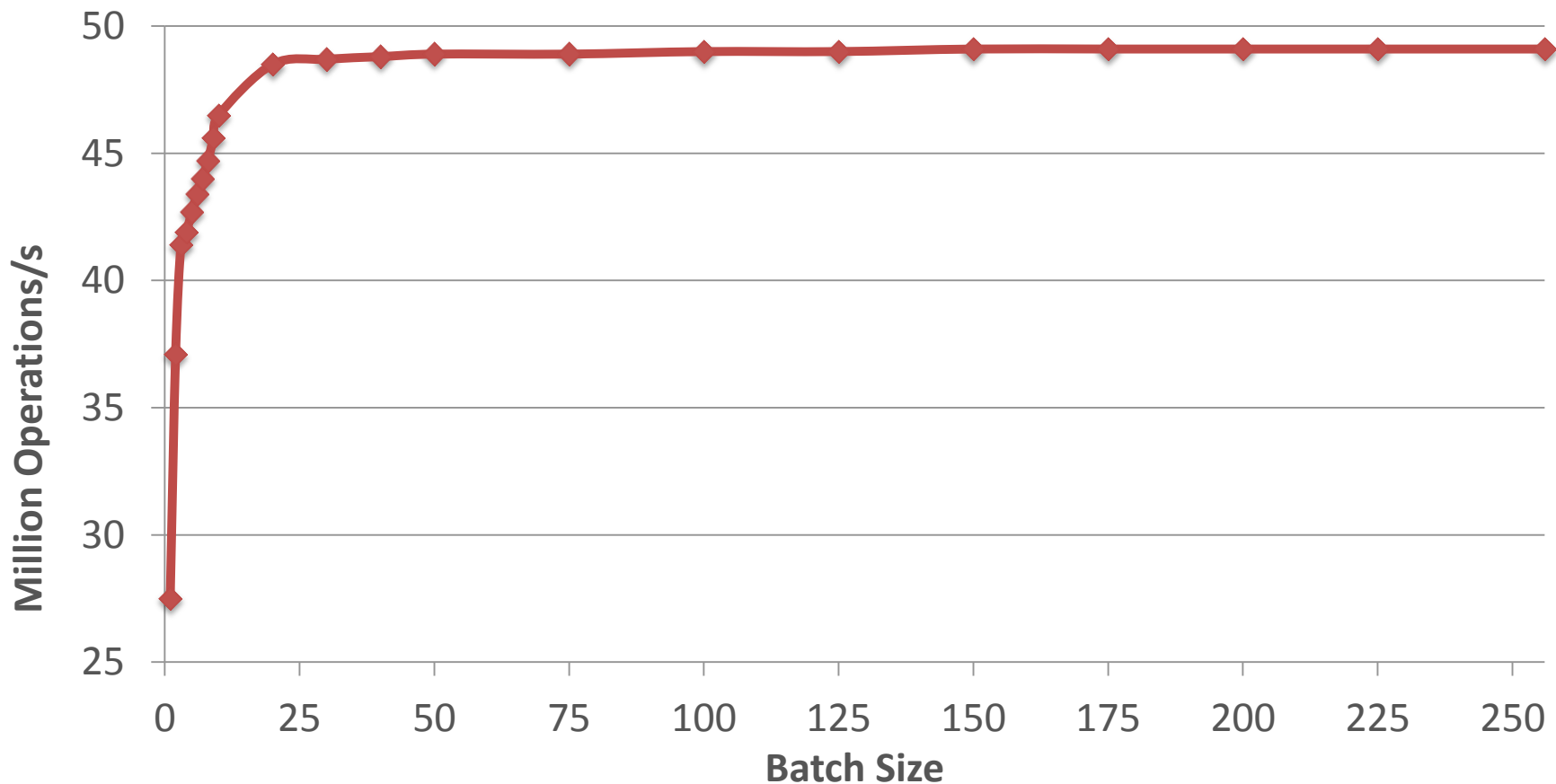
Database **Technology**
Group

- *Intel i7-2600 (4 cores @3.4GHz, 2 memory channels, 16GB DDR3 @1333MHz)*
- *20M keys (32bit), 256-batch size*

- *Intel i7-2600 (4 cores @3.4GHz, 2 memory channels, 16GB DDR3 @1333MHz)*
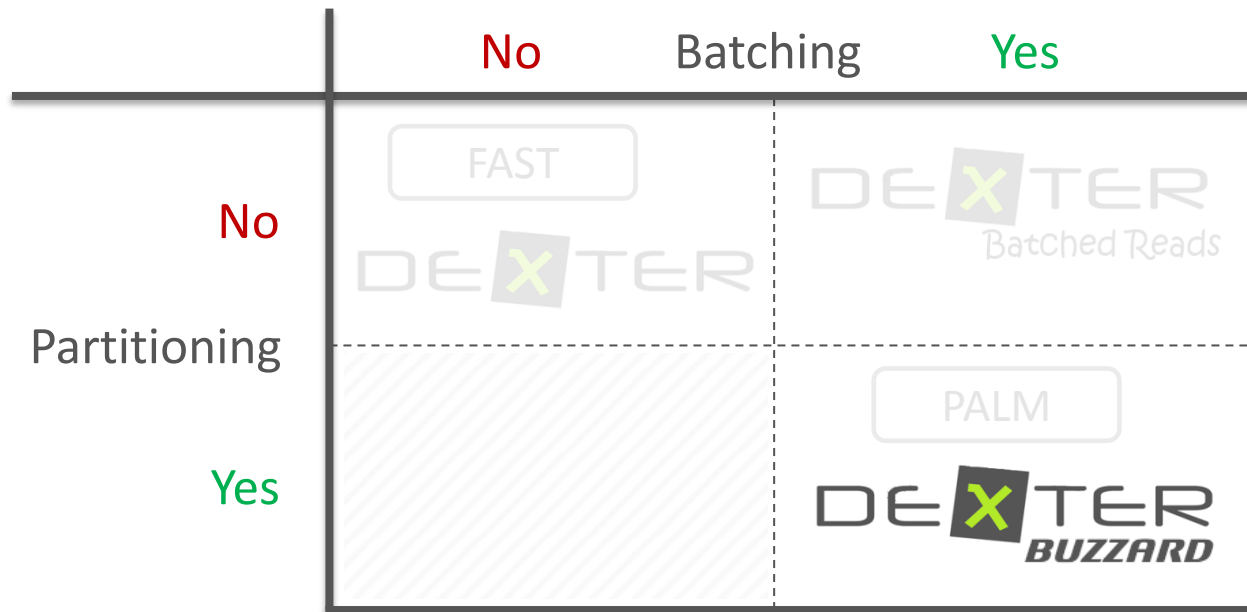- *20M keys (32bit), 256-batch size*

- *Intel i7-2600 (4 cores @3.4GHz, 2 memory channels, 16GB DDR3 @1333MHz)*
- *20M keys (32bit), dynamic batch size*

■ *Two important dimensions for parallel in-memory indexing*

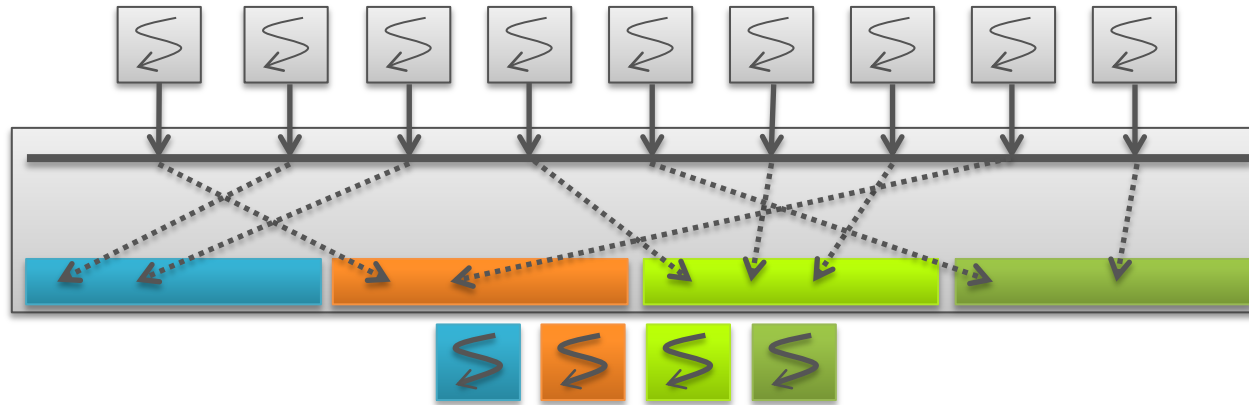|  | No    Batching    Yes |  |
|---|---|---|
| No |  FAST    DEXTER | DEXTER Batched Reads |
| Partitioning | | |
| Yes |  | PALM    DEXTER BUZZARD |

➕ Direct Processing

➖ Synchronization

➕ Independent Data
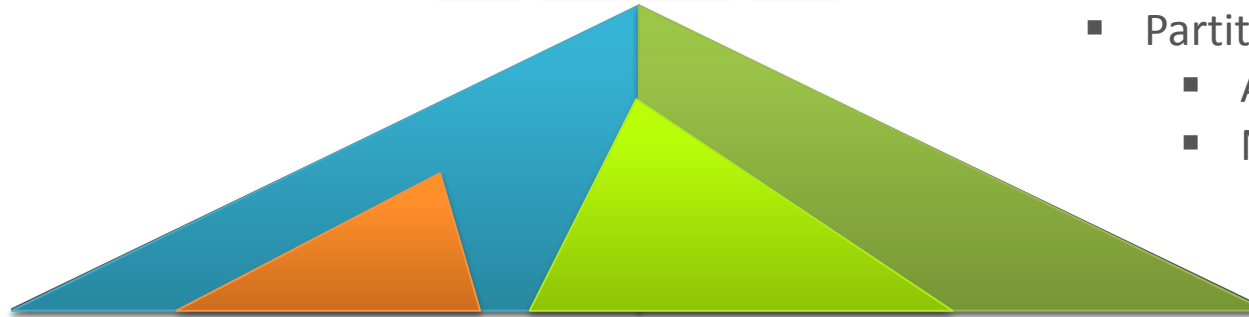
➖ Indirect Processing

➕ Low latency        ➕ High Throughput
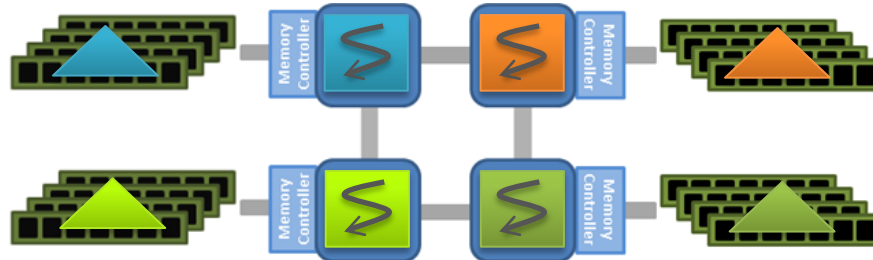
➖ Increased Latency

Incoming Requests

- Thread Mediation Layer
  - Decoupling

- Partitioned Prefix Tree
  - Adaptive partitioning
  - No synchronization needed

- Memory Mapping
  - Map tree partition to local memory on NUMA systems

- Exploit prefix trees for efficient analytical queries
- Database operations on prefix trees

high point query throughput

high update rates

*Direct Query Processing*

*Core Indexing*

- *Idea: Utilize functionality of the underlying hardware and OS to get fast copies of prefix trees for destructive DB operators*

- *Original Prefix Tree*

- *mmapped shared memory*
    - References old VMM
    - Integrates into address space



<< ms

VMM

Memory

VMM

ms

seconds

- *fork*
    - Copies VMM
    - Complex to control

- *memcpy*
    - Copies physical memory
    - Fully independent data

- *Idea: Utilize functionality of the underlying hardware and OS to get fast copies of prefix trees for destructive DB operators*

- *Original Prefix Tree*

- *mmapped shared memory*
    - References old VMM
    - Integrates into address space



<< ms

VMM

Memory

VMM

ms

seconds

- *fork*
    - Copies VMM
    - Complex to control

- *memcpy*
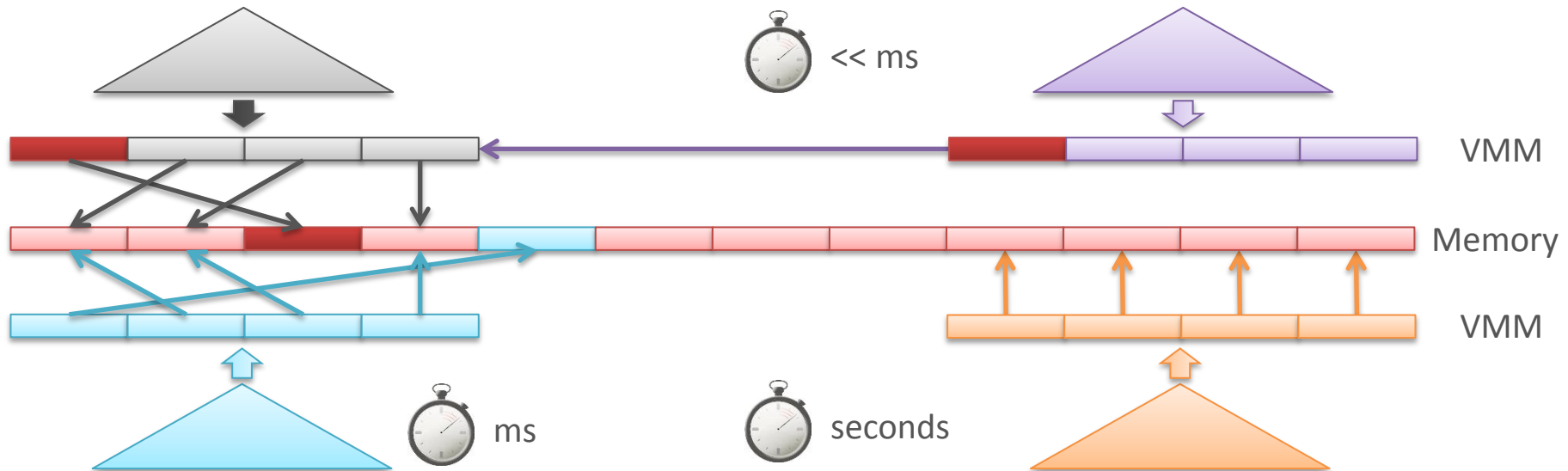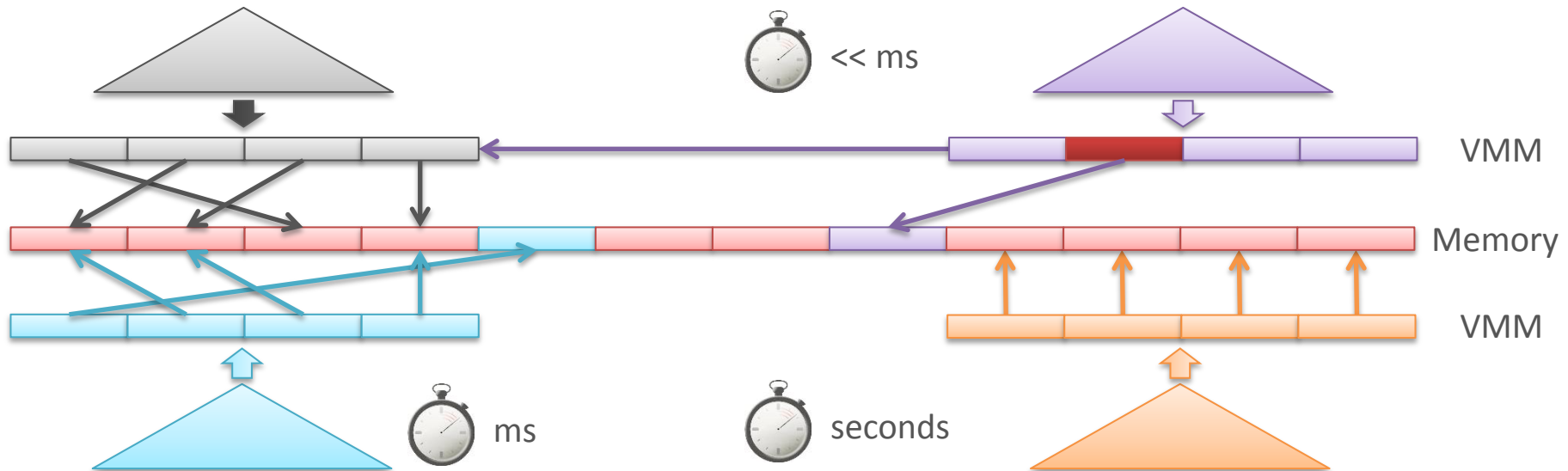    - Copies physical memory
    - Fully independent data

- *Idea: Utilize functionality of the underlying hardware and OS to get fast copies of prefix trees for destructive DB operators*

- *Original Prefix Tree*

- *mmapped shared memory*
  - References old VMM
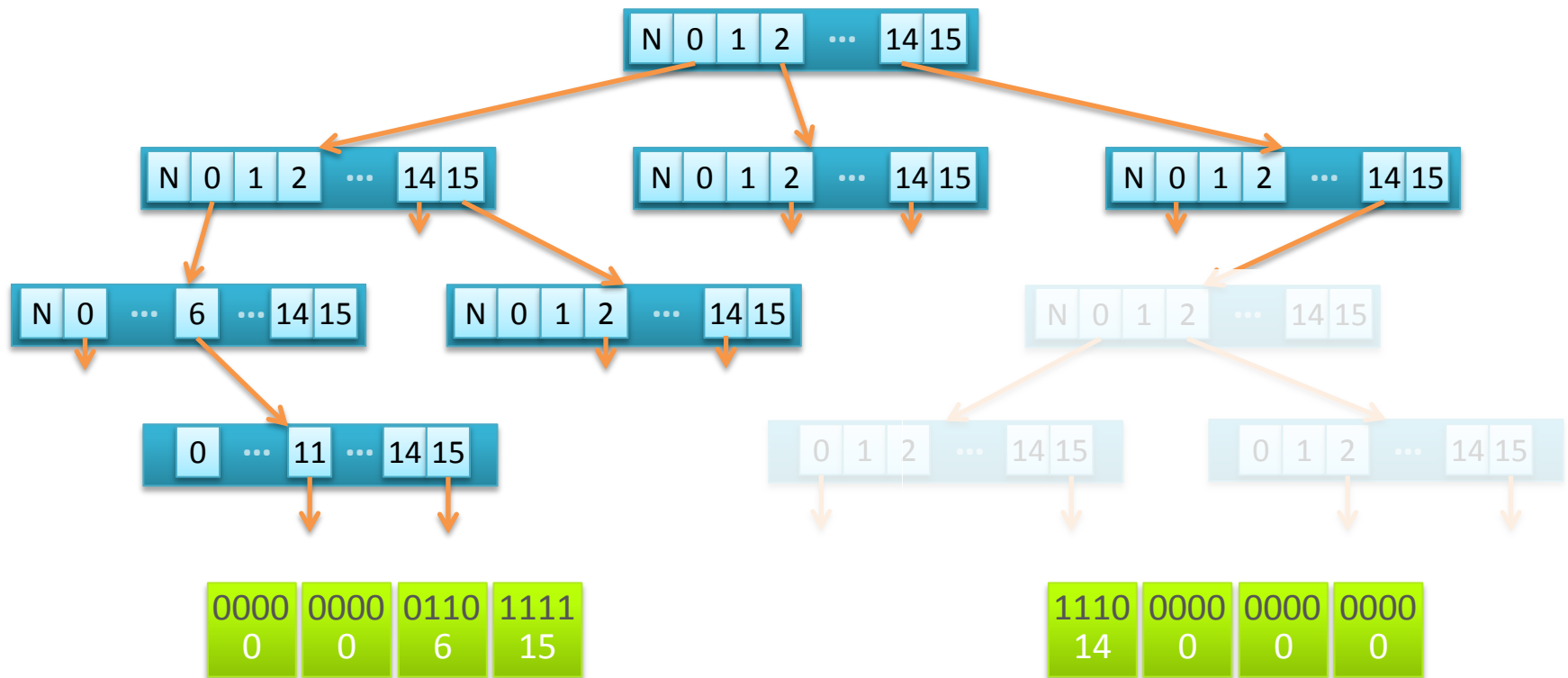  - Integrates into address space

<< ms

VMM

Memory

VMM

ms

seconds

- *fork*
  - Copies VMM
  - Complex to control

- *memcpy*
  - Copies physical memory
  - Fully independent data

- Example for a destructive unary DB operator
- Range selection:  $\sigma_{111 \le A \le 57344}$



- Prune beyond left key path
- Prune beyond right key path

❖ INTRODUCTION AND TRENDS

❖ DEXTER PROJECT

❖ CORE INDEXING

❖ DIRECT QUERY PROCESSING

❖ CONCLUSIONS

### Requirements from Trends

- *Mixed workloads, high update rates, and high responsiveness*
- *Increased parallelism*

### Summary

- *Generalized prefix tree with balanced read/write performance*
- *Direct Query Processing on prefix trees*
- *Main concepts: exploit prefix tree structure for pruning*
- *Promising results and plenty of open work*

### Conclusions

- *Column stores do not fit the needs of Operational BI and Advanced Analytics*
- *Read/write balanced index structures have a high potential for filling this gap*

# Thomas Kissinger

## *DEXTER*
*Parallel In-Memory Indexing and*
*Direct Query Processing on Prefix Trees*

*http://wwwdb.inf.tu-dresden.de/dexter*