



Adaptive Query Processing on Prefix Trees

Wolfgang Lehner

Fachgruppentreffen, 22.11.2012 – TU München

> Challenges for Database Systems

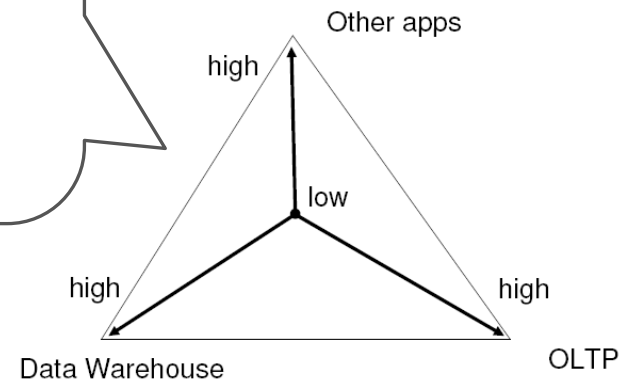


Extreme data



“Three things are important in the database world: performance, performance and performance.”
Bruce Lindsey

The DBMS Landscape – Performance Needs



Extreme performance



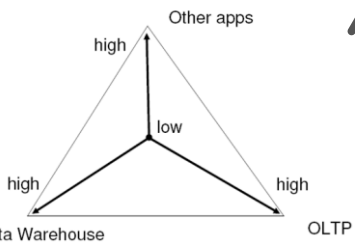
> Challenges for Database Systems



Flexibility in Database Systems

- + during deployment time (schema definition)
- during database lifetime (schema evolution)
- during query runtime (scheduling, ...)

Extreme Performance



Extreme Flexibility



Extreme Data

> Flexibility from 10.000 feet



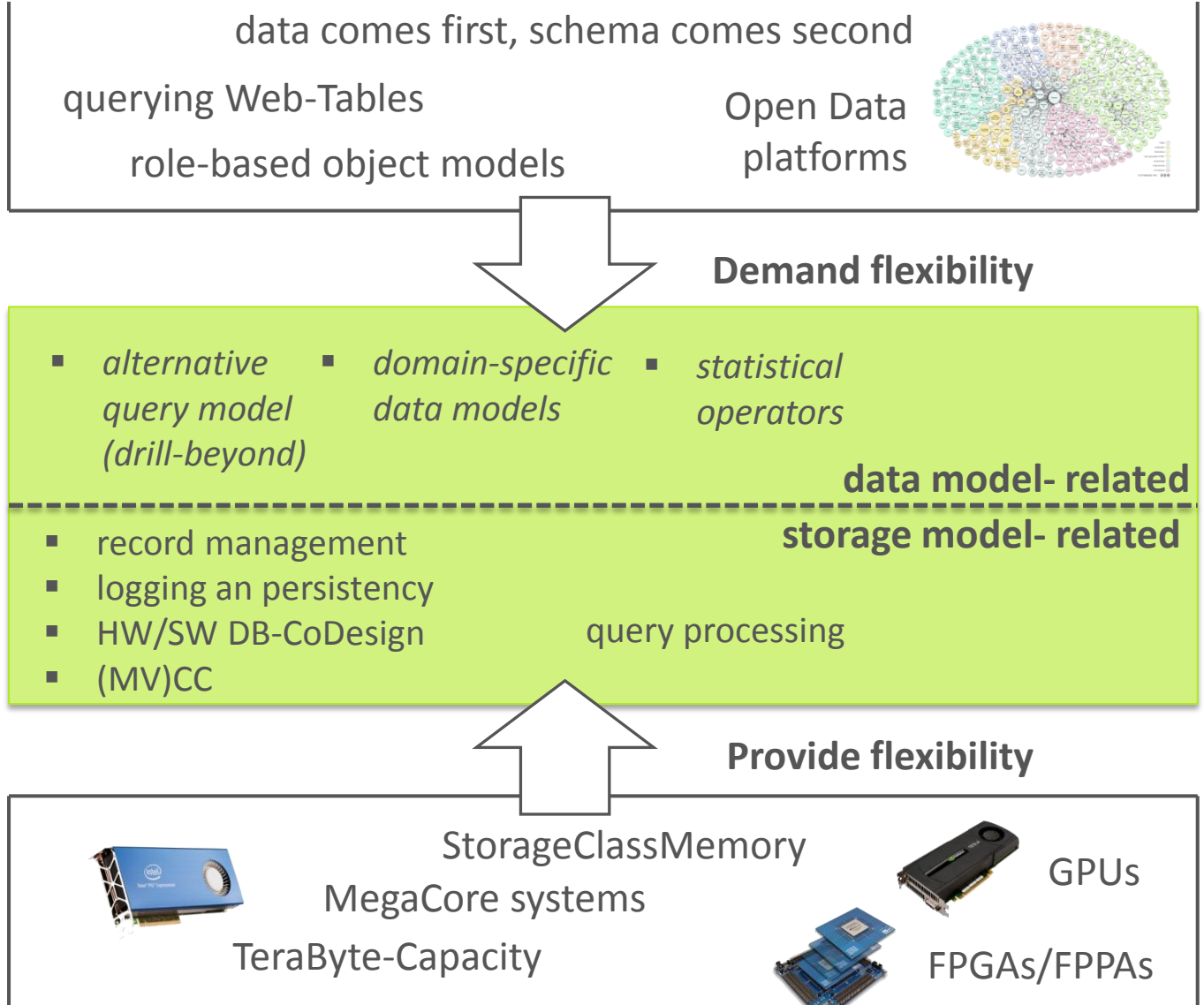
applications

↑

Flexibility is cross-cutting

↓

operating system & hardware



> Example 1: Data-driven Schema Design



1998



2002



2006



2008



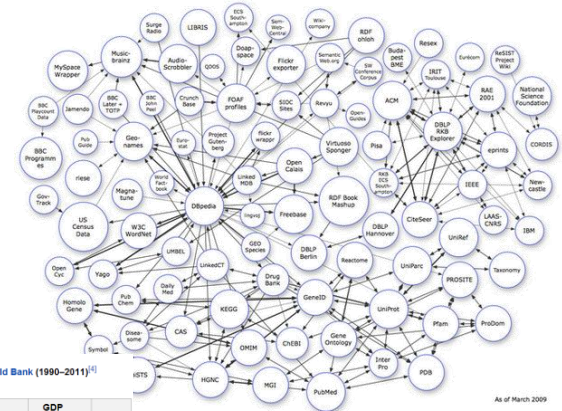
2012

GSM	GSM TriBand	GSM, GPRS, Bluetooth, FM	GSM, UMTS, Bluetooth, USB, WiFi	GSM, UMTS, HSPA+, Bluetooth, FM, USB, WiFi, aGPS, HDMI
Dimensions/Weight	Dimensions/Weight/ Form Factor	Dimensions/Weight/ Form Factor	Dimensions/Weight	Dimensions/Weight
DotMatrix Display	DotMatrix Display	2.0", QVGA, Color	3.5", 320×480, Touch	4.3", 800×480, Touch
Standby/Talk Time	Standby/Talk Time	Standby/Talk Time	Battery Capacity	Battery Capacity
#Names, #SMS	#Names, #SMS	Memory/Card Slot	Memory	Memory/Storage/Car rd Slot
	Polyphonic Ringtones	MP3 Ringtones/MP3 Player/Browser	Browser/Exchange	Browser/Exchange
		3MP Camera	3MP Camera	Front/Rear Camera/HD Video
			CPU, OS	CPU, GPU, OS



Web Data: A Multitude of Data Models and Sources

OAuth
JSON
REST
Basic-Auth
DELETE



List by the United Nations (2010)^[2]

List by the United Nations (2010)^[2]

List by the World Bank (1990-2011)^[4]

Rank	Country/Region	2010 GDP (millions of US\$)	Rank	Country/Region	GDP (millions of US\$)	Rank	Country/Region	GDP (millions of US\$)	Year
	World	62,633,783		World	69,899,225 ^[5]		World	69,983,693 ^[12]	2011
1	United States	14,447,100	1	United States	15,075,675	1	United States	15,094,000	2011
2	China	5,739,358	2	China	7,298,147 ^[2]	2	China	7,318,499 ^[2]	2011
3	Japan	5,458,873	3	Japan	5,866,540	3	Japan	5,867,154	2011
4	Germany	3,280,334	4	Germany	3,607,364	4	Germany	3,570,556	2011
5	France	2,559,850	5	France	2,778,085	5	France	2,773,032 ^[4]	2011
6	United Kingdom	2,253,552	6	Brazil	2,492,907	6	Brazil	2,476,652	2011
7	Brazil	2,088,966	7	United Kingdom	2,431,310	7	United Kingdom	2,431,589	2011
8	Italy	2,051,290	8	Italy	2,198,732	8	Italy	2,194,750	2011
9	India	1,722,328	9	Russia	1,850,401	9	Russia	1,857,770	2011
10	Canada	1,577,040	10	India	1,826,811	10	India	1,847,982	2011
11	Russia	1,479,823	11	Canada	1,738,954	11	Canada	1,736,051	2011
12	Spain	1,407,322	12	Australia	1,486,914	12	Spain	1,490,810	2011
13	Australia	1,271,945	13	Spain	1,479,560	13	Australia	1,371,764	2011
14	Mexico	1,032,224	14	Mexico	1,153,958	14	Mexico	1,155,316	2011
15	South Korea	1,014,369							

- Example 1: Queries spanning millions of data sources with as many schemata
- Example 2: Augment local datasets with (external) third-party information



select ... from cust,
 where nations.GDP > \$1B

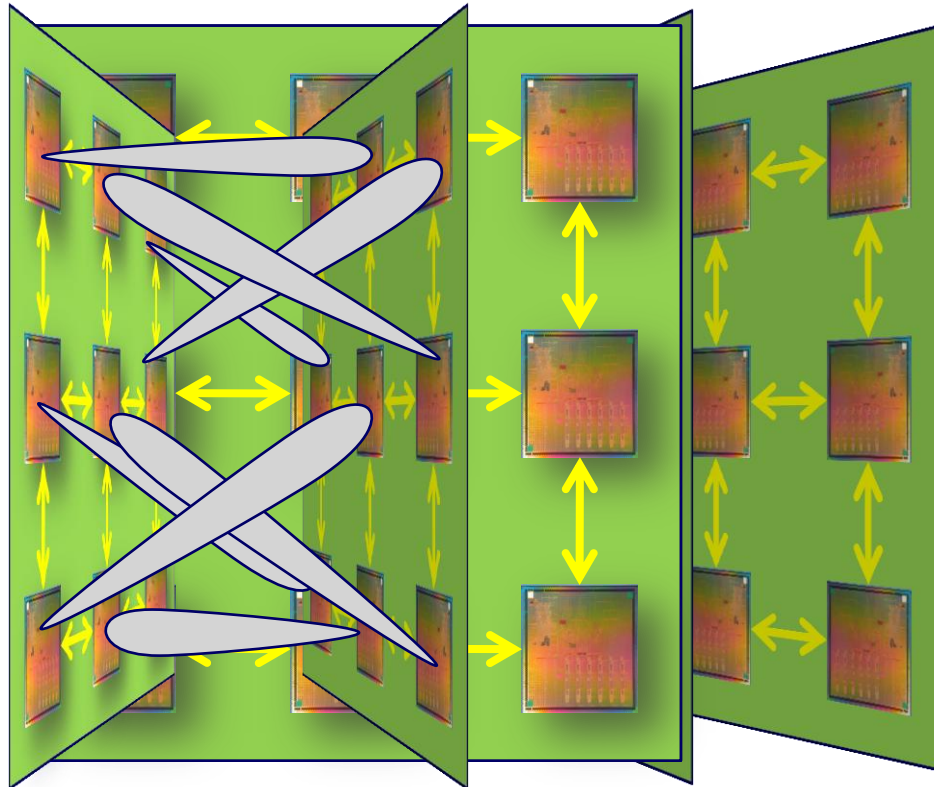


> Example 3: Energy-Adaptive High-Speed Computing Platform



Deutsche
Forschungsgemeinschaft
DFG

SFB912: HAEC – Highly Adaptive Energy Efficient Computing



Optical Interconnect

- adaptive analog/digital circuits for e/o transceiver
- embedded polymer waveguide
- packaging technologies (e.g. 3D stacking of Si/III-V hybrids)
- 90° coupling of laser

Radio Interconnect

- on-chip/on-package antenna arrays
- analog/digital beamsteering and interference minimization
- 100Gb/s
- 100-300GHz channel
- 3D routing & flow management





3

Composed Operators

- Construction of composed operators, like multi-way-select-join-group-sort
- Avoids Intermediate Materialization

1

Intermediate Indexed Tables

- Operators exchange clustered indexes

2

Cooperative Operators

- Operators optimize their output for the next operator

CIDR2013

Thomas Kissinger, Benjamin Schlegel, Dirk Habich, Wolfgang Lehner
 Database Technology Group
 Technische Universität Dresden
 01062 Dresden, Germany
 (firstname.lastname)@tu-dresden.de

ABSTRACT

Modern database systems have to process large amounts of data and should provide results with low latency at the same time. To achieve this, data is nowadays typically held completely in main memory, to benefit of its high bandwidth and low access latency that could never be reached with disks. Current in-memory databases are usually column-stores that exchange columns or vectors between operators and suffer from a high tuple reconstruction overhead. In this paper, we present the *indexed table-at-a-time* processing model that makes indexes the first-class citizen of the database system. The processing model comprises the concepts of intermediate indexed tables and cooperative operators, which make indexes the common data exchange format between plan operators. To keep the intermediate index materialization costs low, we employ optimized prefix trees that offers a balanced read/write performance. This processing model allows the construction composed operators like the multi-way-select-join-group. Such operators speed up the processing of complex OLAP queries so that our approach outperforms state-of-the-art in-memory databases.

1. INTRODUCTION

Processing complex OLAP queries with low latencies as well as high throughput on ever-growing, huge amounts of data is still a major challenging issue for modern database systems. With the general availability of high main memory capacity, in-memory database systems become more and

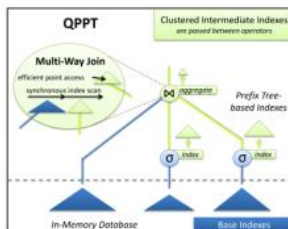
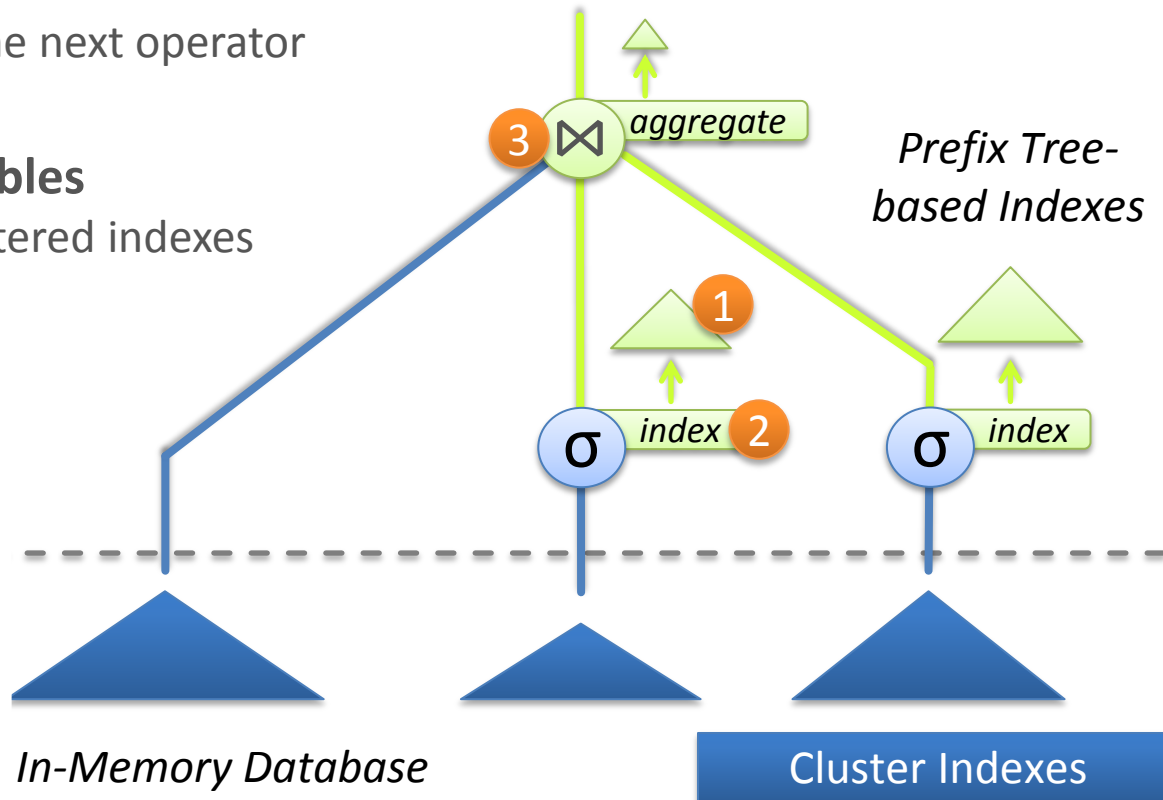
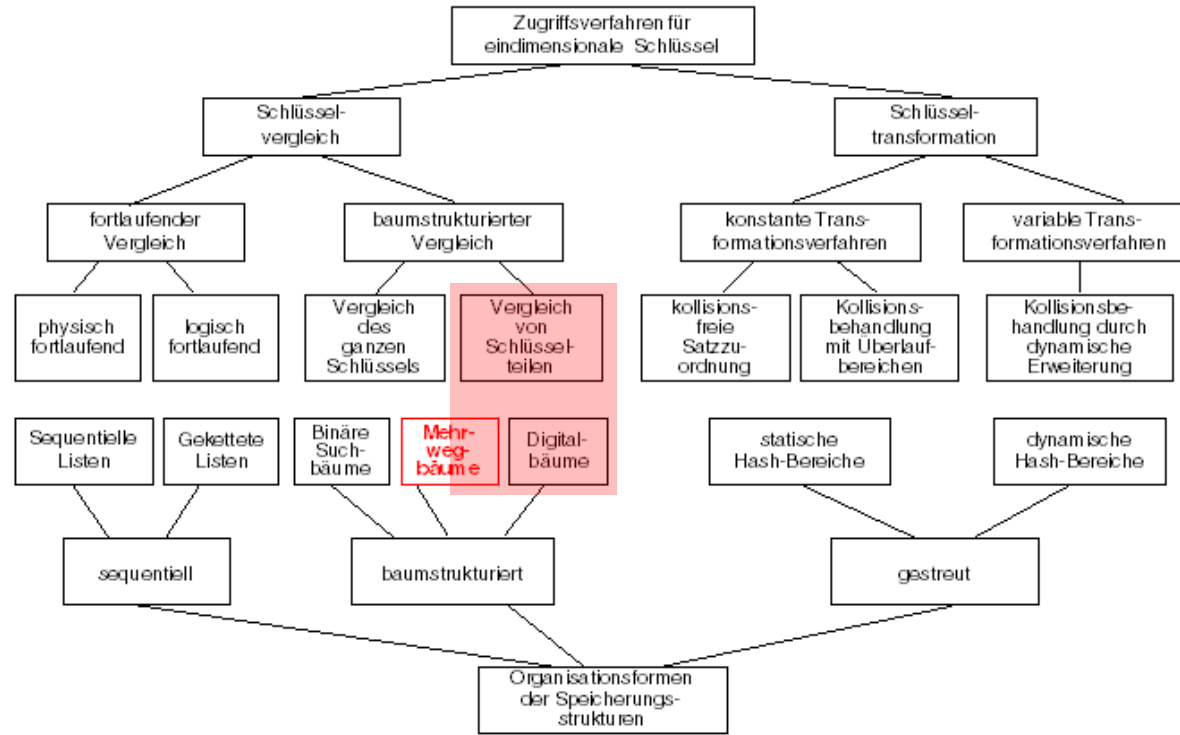


Figure 1: Overview of QPPT's indexed table-at-a-time Processing Model.

processing (e.g., massive virtual function calls) and are thus more suitable for modern in-memory column-stores. However, the logical unit of a tuple as present in row-stores is decomposed in columns resulting in an expensive tuple reconstruction overhead. The complexity of the tuple reconstruction procedure grows with an increasing number of





Prefix Trees

- Unbalanced
- + Balanced read/write performance
- + Main memory-optimized
- + Efficient parallel access

B+-Trees

- Balanced
- Low update performance
- Disk-optimized
- Low scalability

> KISS-Tree Overview



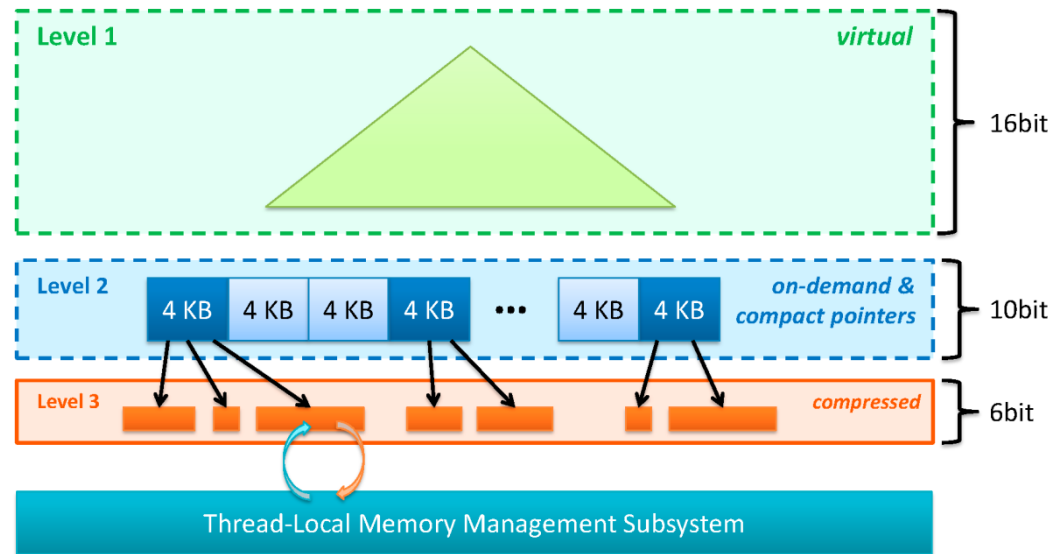
Properties

- Specialized version for 32bit keys
- Latch-free updates
- Order-preserving
- 2-3 memory accesses per key

→ Comparable fast to reported order-preserving in-memory indexes for read access

→ BUT:
High update performance

- Heterogeneous in-memory index structure
- Combination of direct and indirect addressing
- Takes advantage of virtual memory management
- Enables different compression mechanisms



KISS-Tree: Smart Latch-Free In-Memory Indexing on Modern Architectures DAEMON 2012

Thomas Kissinger, Benjamin Schlegel, Dirk Habich, Wolfgang Lehner
Database Technology Group
Technische Universität Dresden
01062 Dresden, Germany
(firstname.lastname)@tu-dresden.de

ABSTRACT

Growing main memory capacities and an increasing number of hardware threads in modern server systems led to fundamental changes in database architectures. Most importantly, query processing is nowadays performed on data that is often completely stored in main memory. Despite of a high main memory scan performance, index structures are still important components, but they have to be designed from scratch to cope with the specific characteristics of main memory and to exploit the high degree of parallelism. Current research mainly focused on adapting block-optimized B+-Trees, but these data structures were designed for secondary memory and involve comprehensive structural maintenance for updates.

In this paper, we present the *KISS-Tree*, a latch-free in-memory index that is optimized for a minimum number of memory accesses and a high number of concurrent updates. More specifically, we aim for the same performance as modern hash-based algorithms but keeping the order-preserving nature of trees. We achieve this by using a prefix tree that incorporates virtual memory management functionality and compression schemes. In our experiments, we evaluate the *KISS-Tree* on different workloads and hardware platforms and compare the results to existing in-memory indexes. The *KISS-Tree* offers the highest reported read performance on current architectures, a balanced read/write performance, and has a low memory footprint.

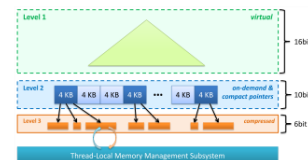
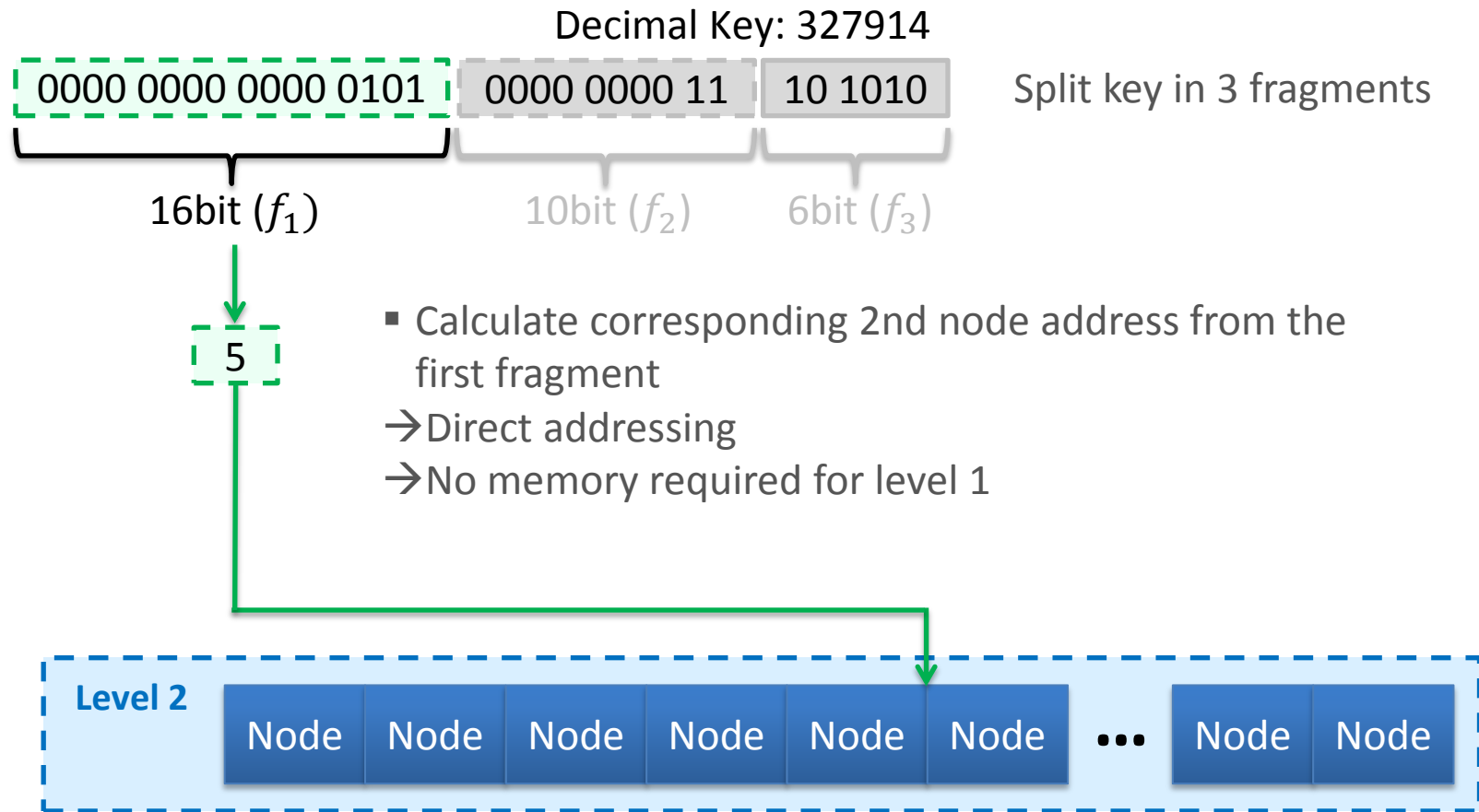


Figure 1: KISS-Tree Overview.

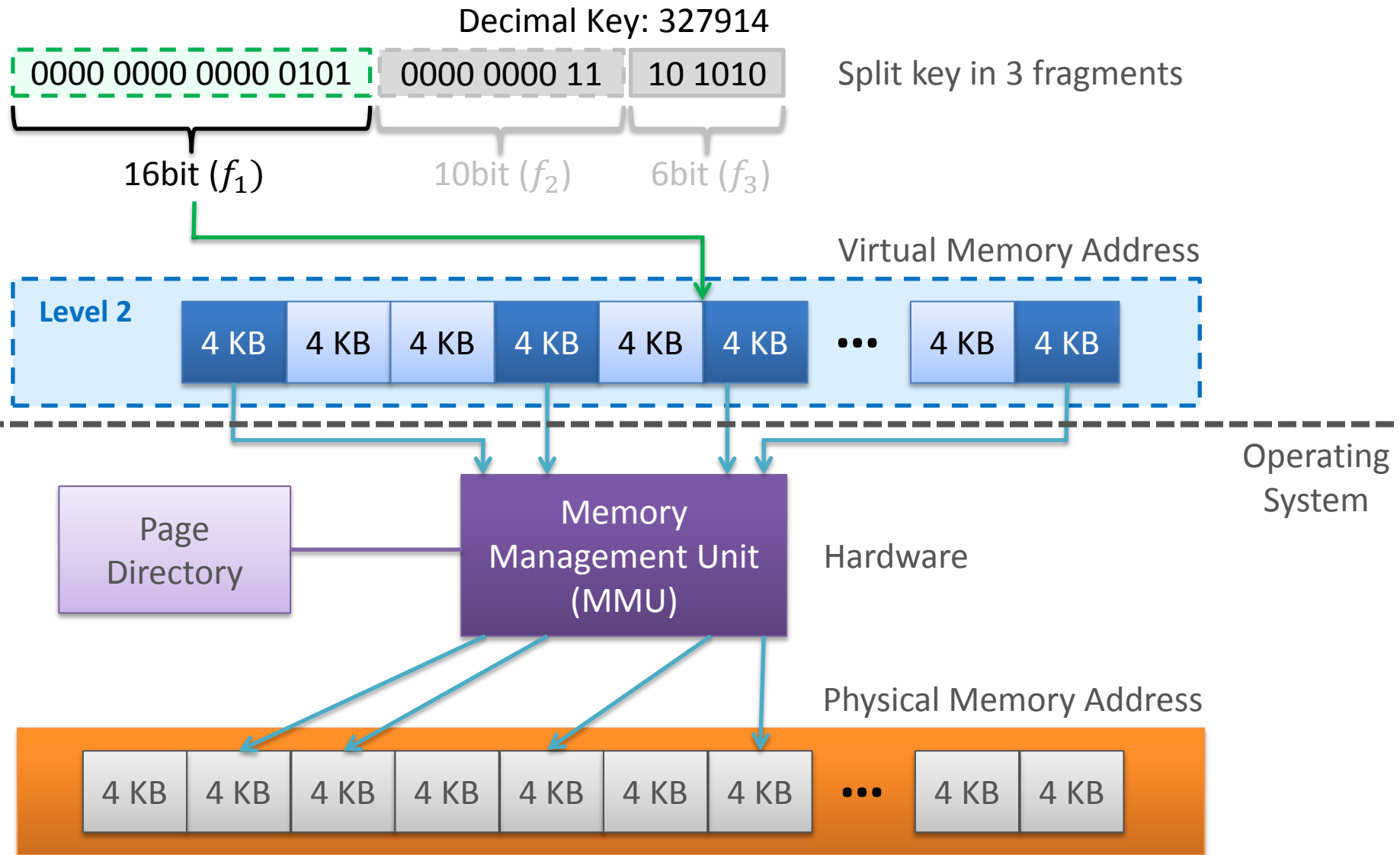
indexes) in-memory and use secondary memory (e.g., disks) only for persistence. While disk block accesses constituted the bottleneck for disk-based systems, modern in-memory databases shift the memory hierarchy closer to the CPU and face the "memory wall" [13] as new bottleneck. Thus, they have to care for CPU data caches, TLBs, main memory accesses and access patterns. This essential change also affects indexes and forces us to design new index structures that are now optimized for these new design goals. Moreover, the movement from disks to main memory dramatically increased data access bandwidth and reduced latency. In combination with the increasing number of cores on modern hardware, parallel processing of operations on index structures imposes a new challenge for us, because the overhead

> Level 1: Virtual Level



Requirement for Level 2: All nodes have to be stored sequentially in memory

> Level 1: Virtual Level

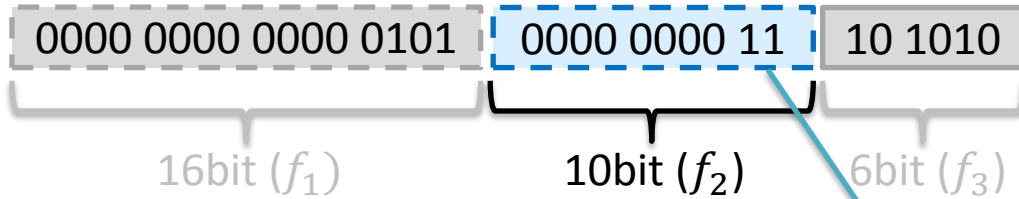


> Level 2: On-demand Level



- Change over to indirect addressing
- 1024 buckets per node containing a compact pointer to the 3rd level node
- 256 MB maximum

Decimal Key: 327914



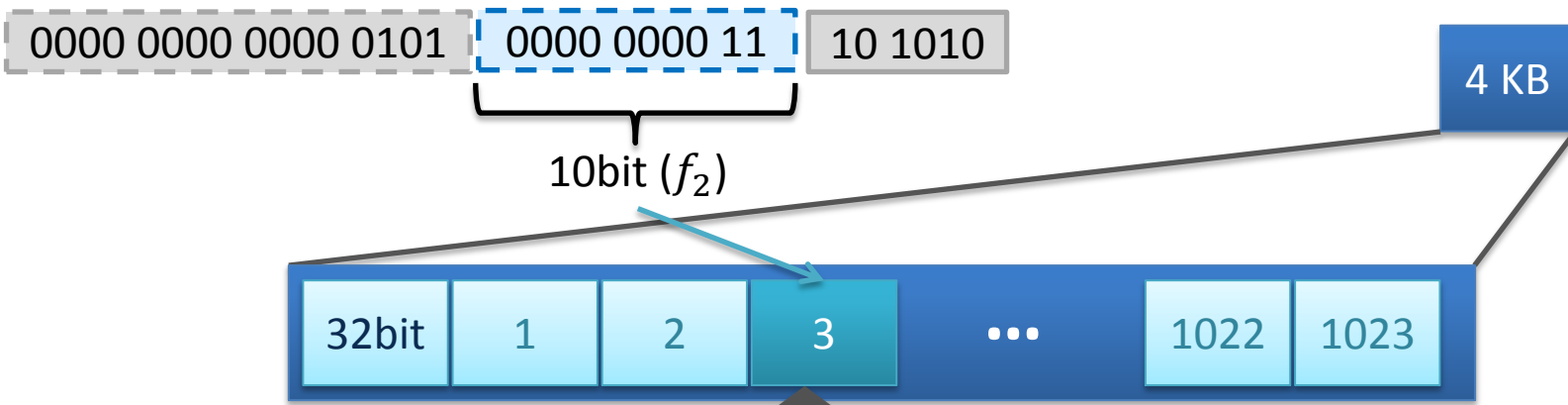
Virtual Memory Address



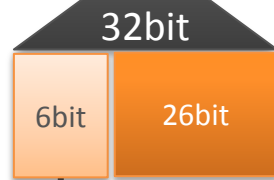
$2^{10} = 1024$ pointer to L3 leaf nodes

$2^{16} + 2^{10} =$ potential L3 leaf nodes of a size between 2^0 and 2^6

> Level 2: On-demand Level



Number of existing elements in L3 leaf node (compression)



Position of L3 leaf node in corresponding memory segment

- Allocation of consecutive virtual memory segments for each of the 64 node sizes possible on the third level
- Each segment consists of a maximum of 2^{26} blocks of the respective node size

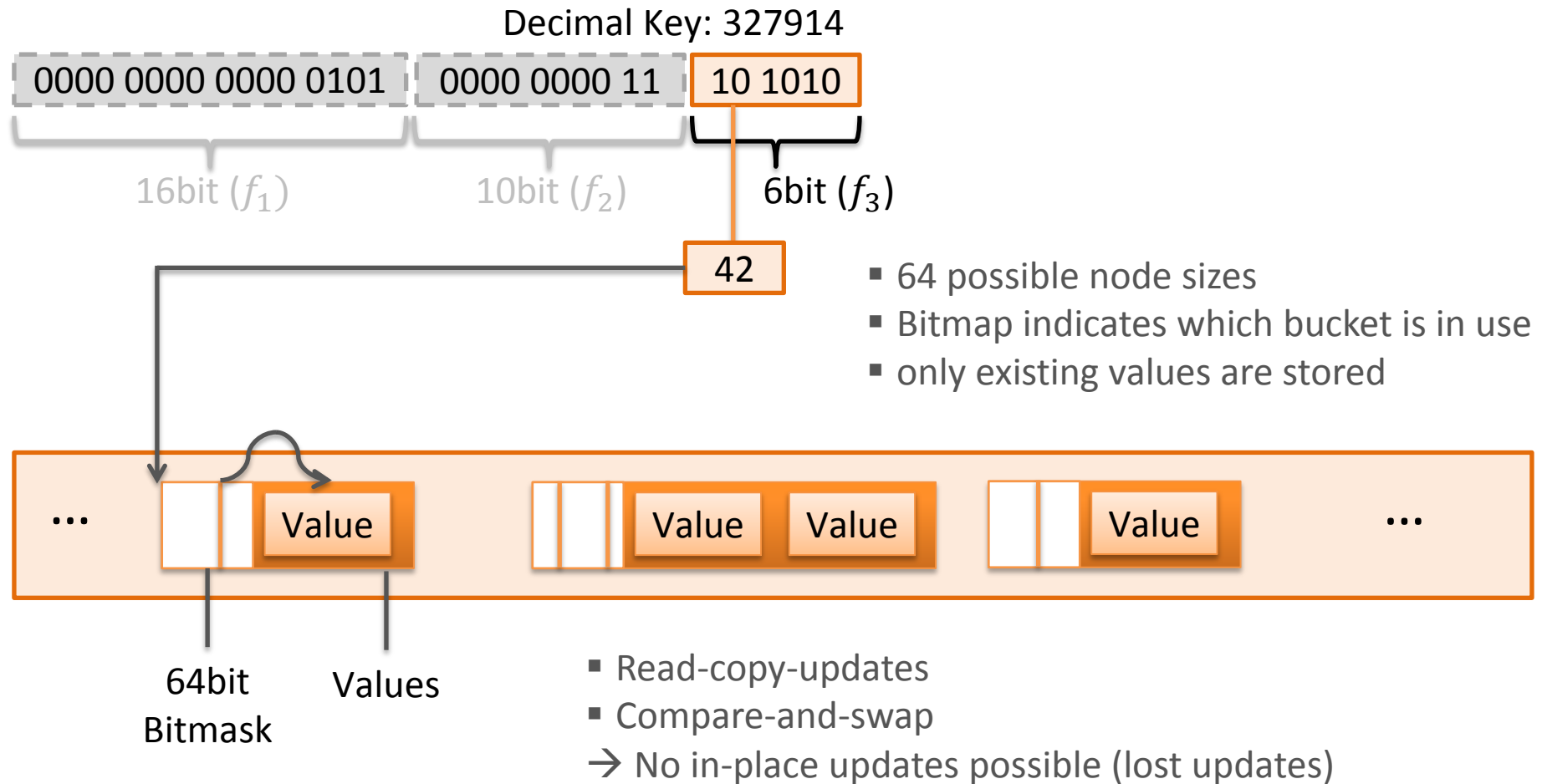


→ 2^{26} blocks of size 1 pre-allocated



→ 2^{26} blocks of node size 64 pre-allocated

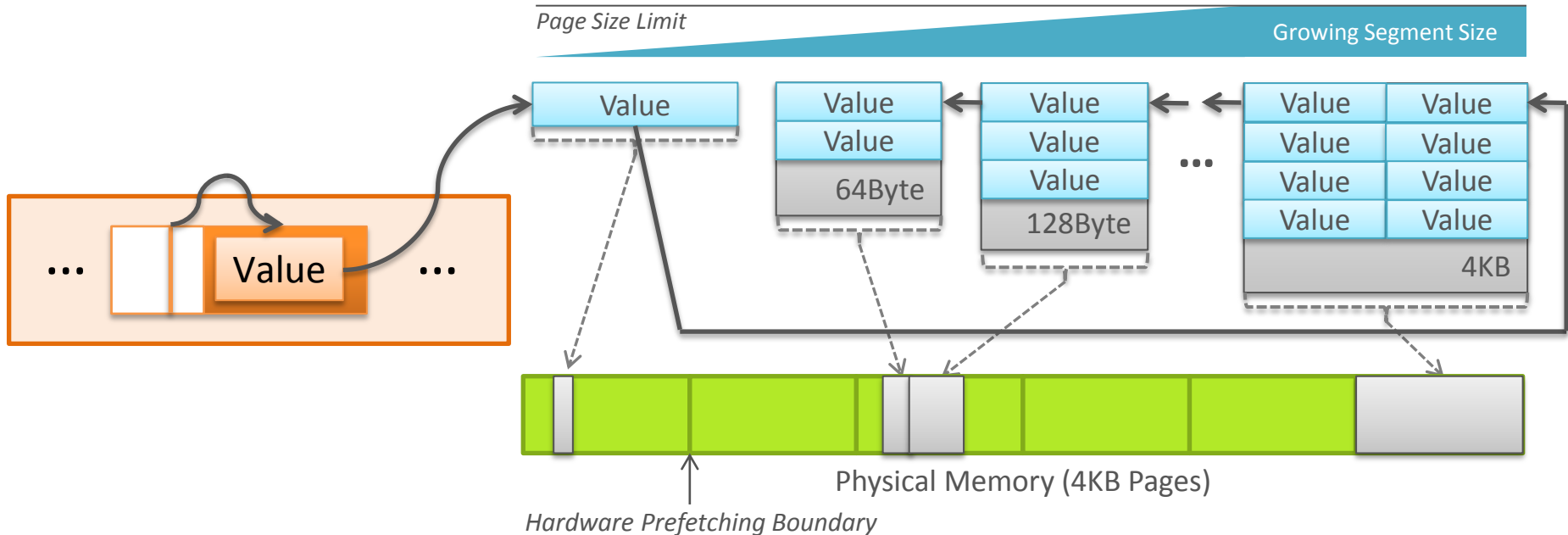
> Level 3: Compression



Thread-Local Memory Management Subsystem



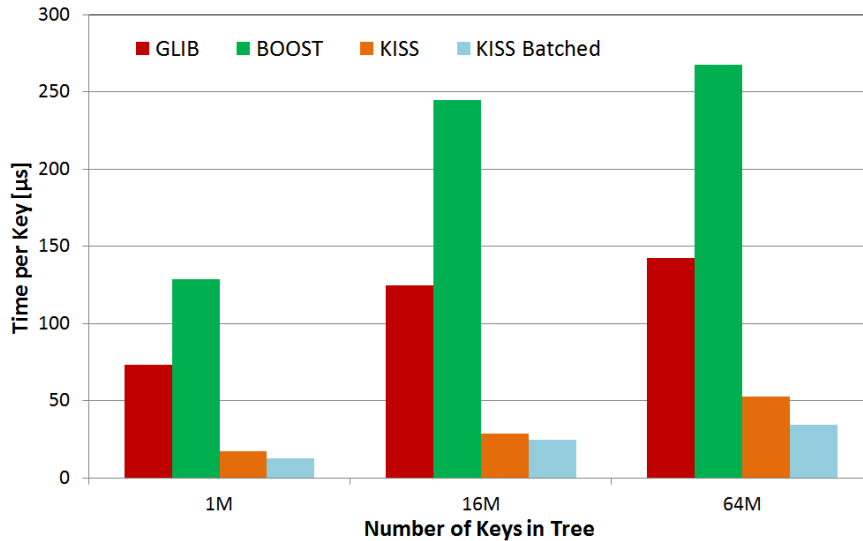
- Efficient duplicate handling necessary for query processing
 - Scanning a linked list results in random memory accesses



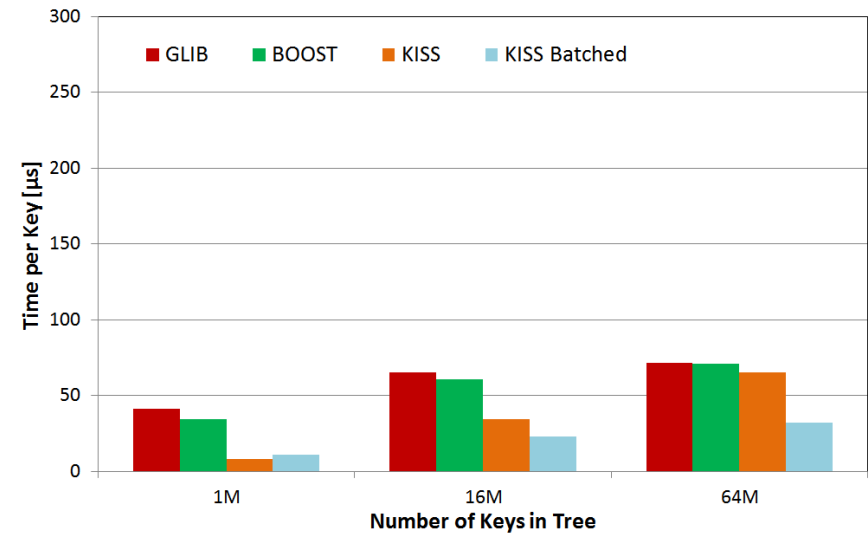
- Page boundaries are a barrier for hardware prefetchers
 - store values sequentially in 4KB blocks
 - blocks grow exponentially until reaching 4KB
 - trade-off between scan performance and memory consumption



Insert/Update Performance



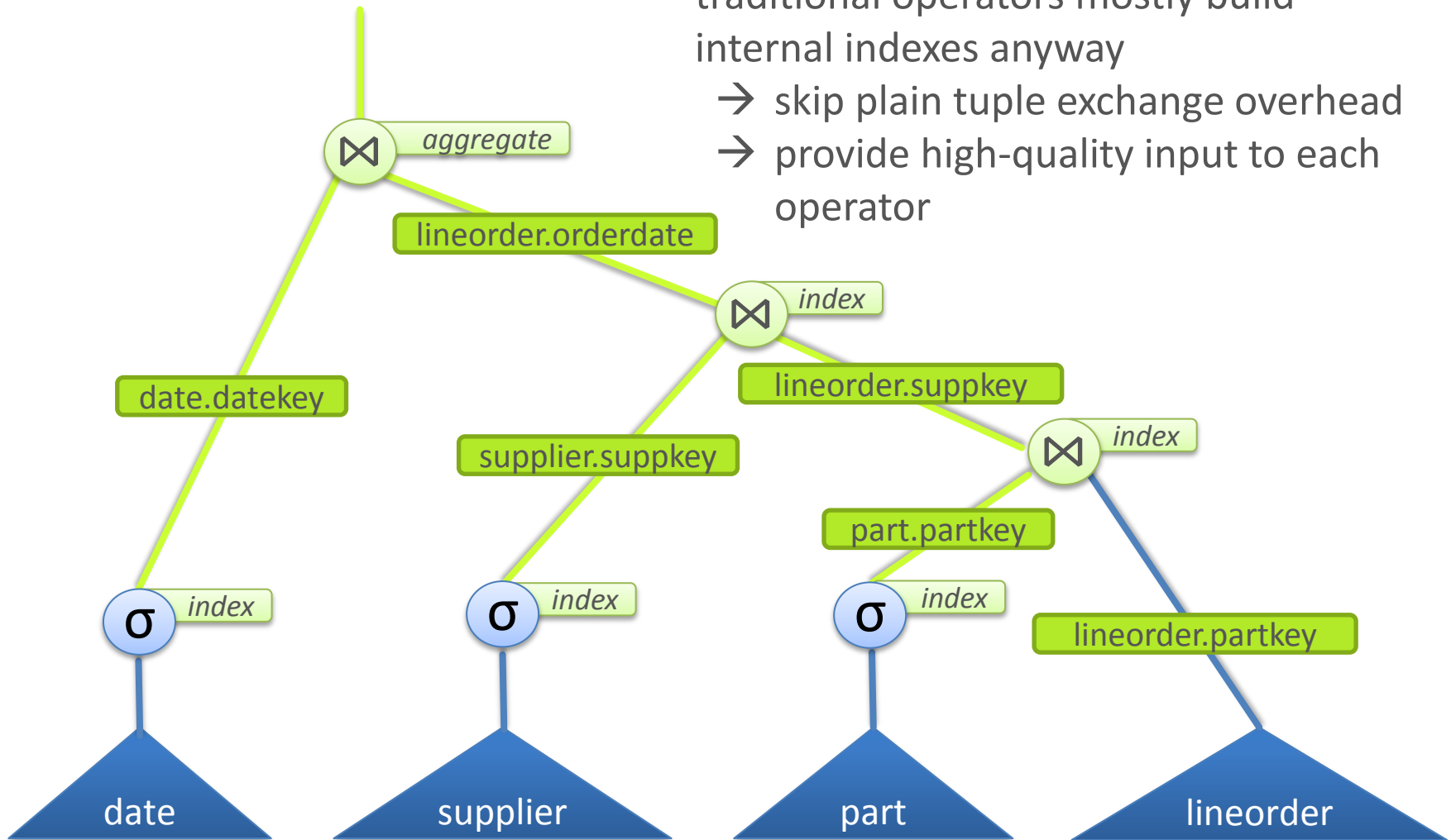
Lookup Performance



- Performance depends on the number of memory accesses per key
 - KISS-Tree: 3 memory accesses per key
 - Latency Hiding with Batch Updates/Lookups
- KISS-Tree performs better than hash table implementations



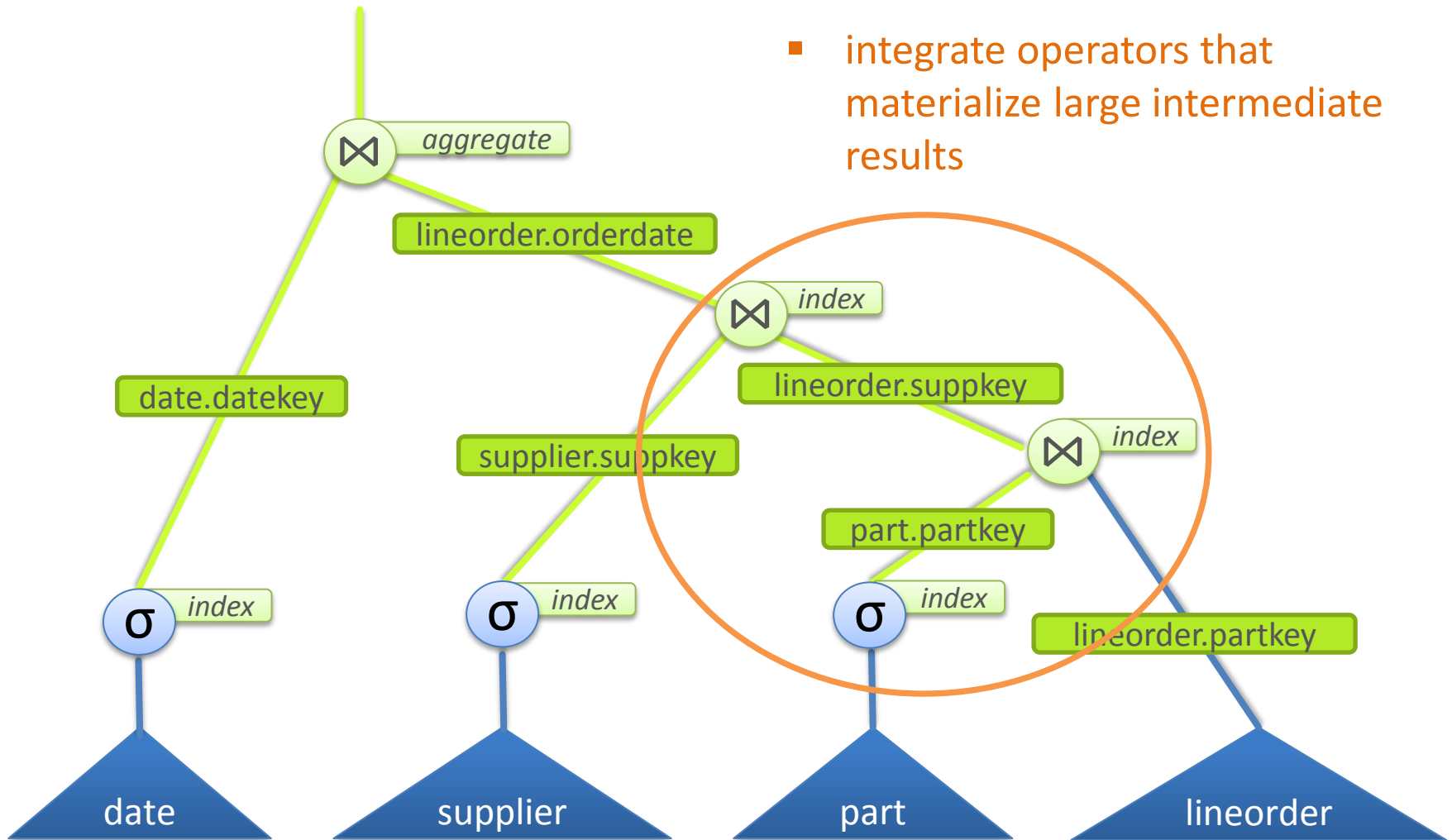
- Each operator adjusts its output to the requirements of the successive operator
 - traditional operators mostly build internal indexes anyway
 - skip plain tuple exchange overhead
 - provide high-quality input to each operator





- Materialization is costly; thus try to completely avoid it!

- integrate operators that materialize large intermediate results

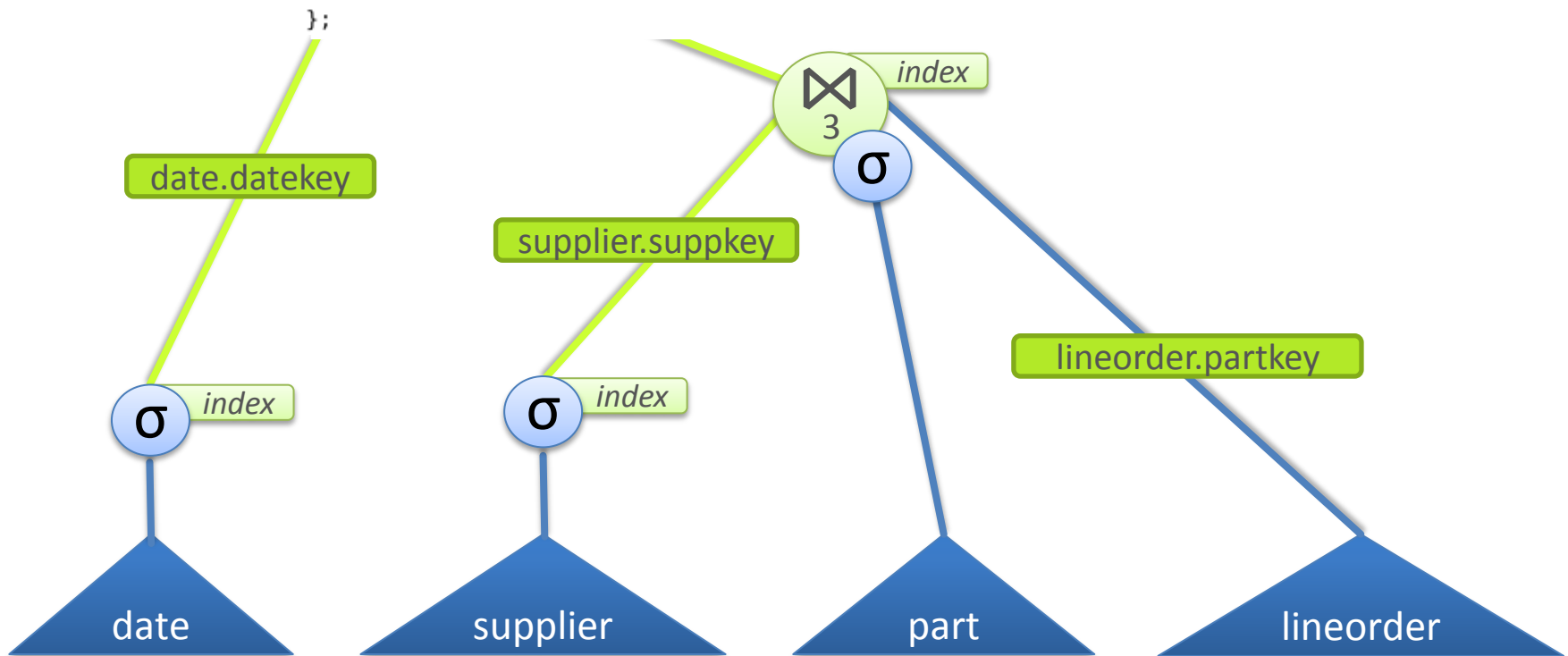




- Materialization is costly; thus try to completely avoid it!

3-way-select-join

```
class FILTERPT4: public PlanOperator {  
    ...  
    void joinWith(int joinnum, int outputnum, u_int64_t** buffer, int bufferpos, bool final);  
    u_int64_t**** joinbuffers;  
    ...  
};
```

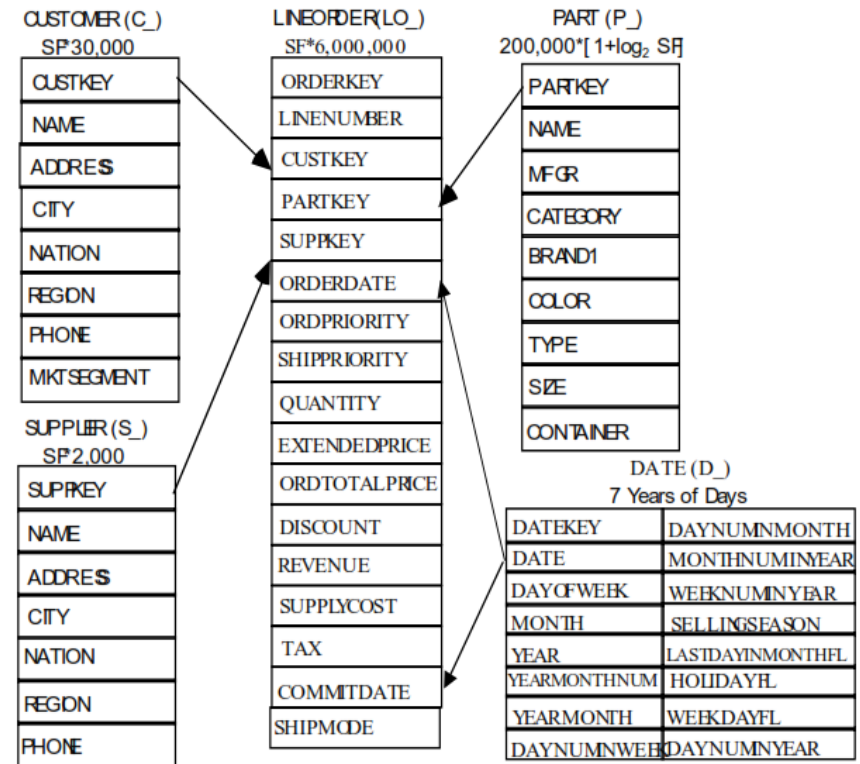


> Evaluation (1)

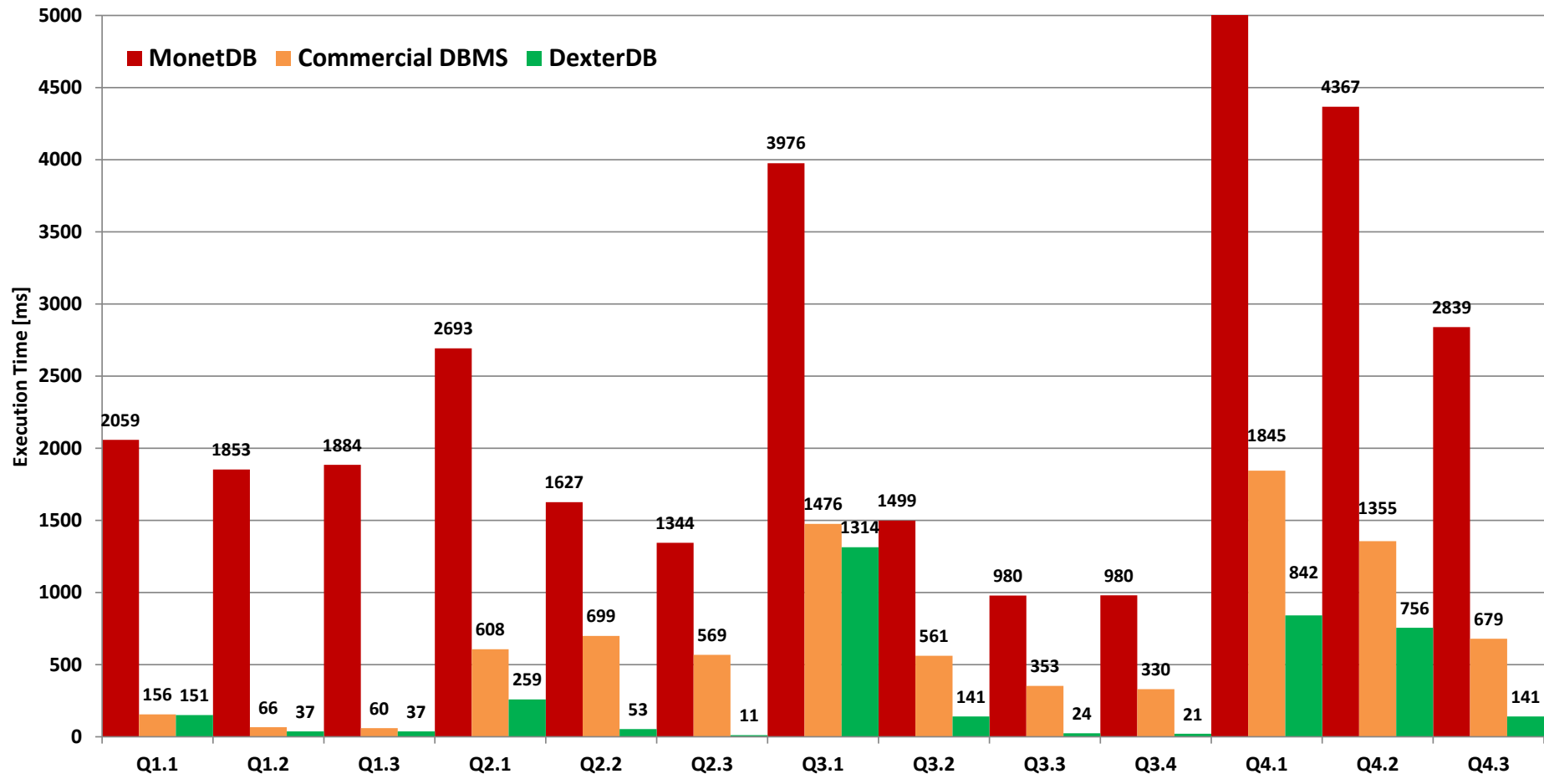


- Star Schema Benchmark (SSB)
 - Scale Factor = 15

- Intel i7-3960X
 - 15 MB LLC
 - 3.3 GHz (Turbo Mode disabled)
- 32 GB Main Memory
- Ubuntu 12.04
- Single-threaded execution

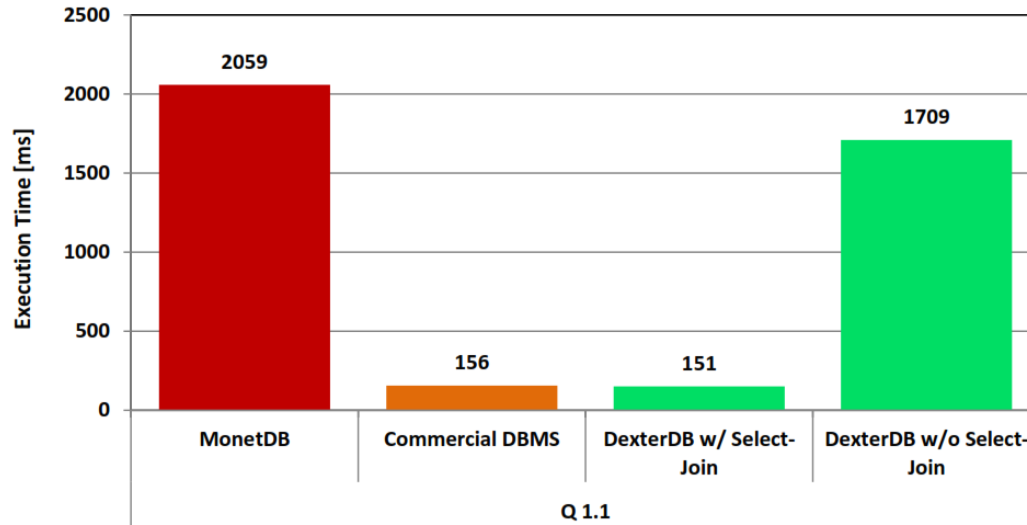


> Evaluation (2)



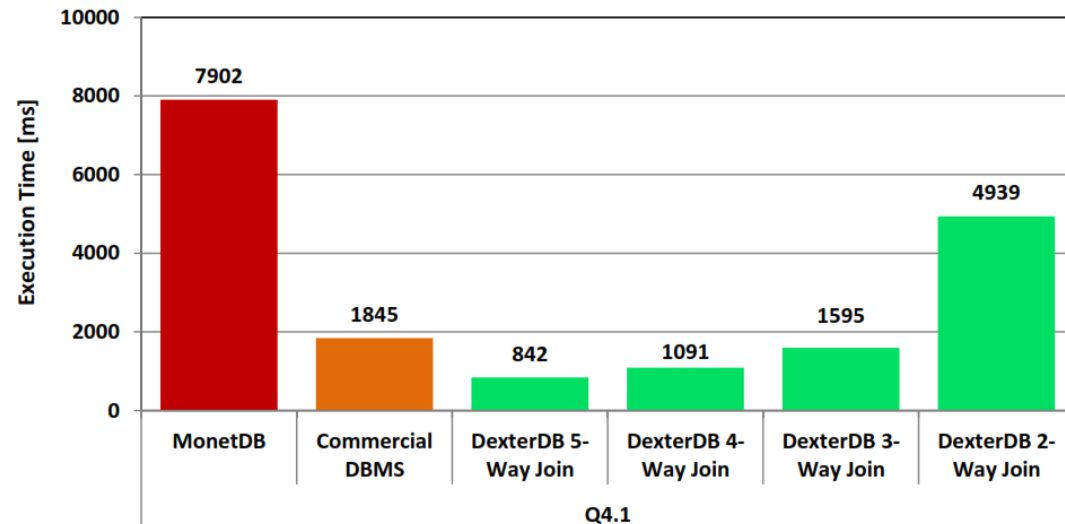
(*) : we are in touch with MonetDB; performance numbers are under investigation

> Evaluation (3)



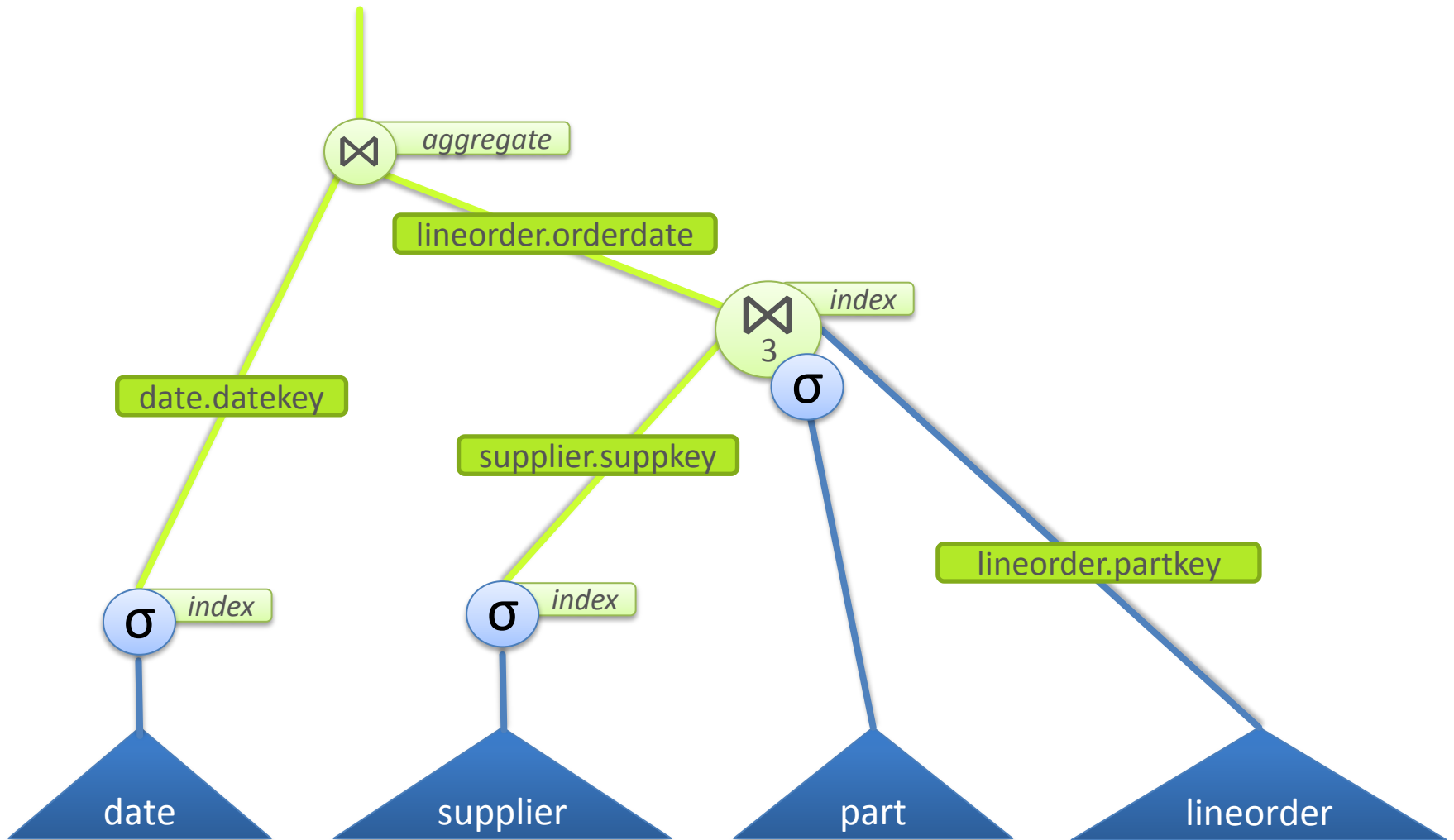
- SSB Query 1.1
- Single Join
- Select-Join operator avoids intermediate result materialization
- Batched KISS-Tree Lookups/Inserts

- SSB Query 4.1
- 4 Joins over 5 Tables
- Multi-way join operator avoids intermediate result materialization

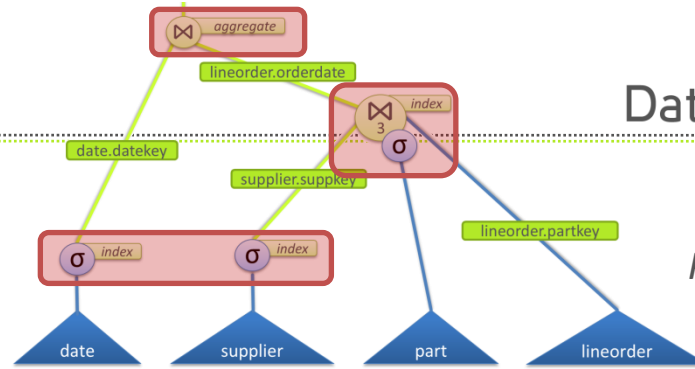




- Flexible Query Execution on Adaptive Hardware



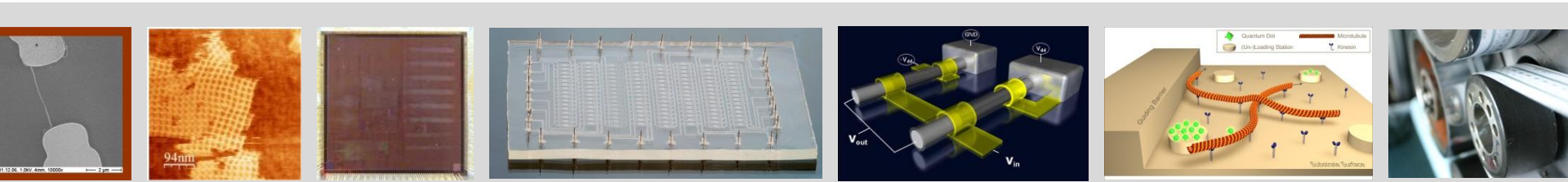
> Adaptive QPPT (2)



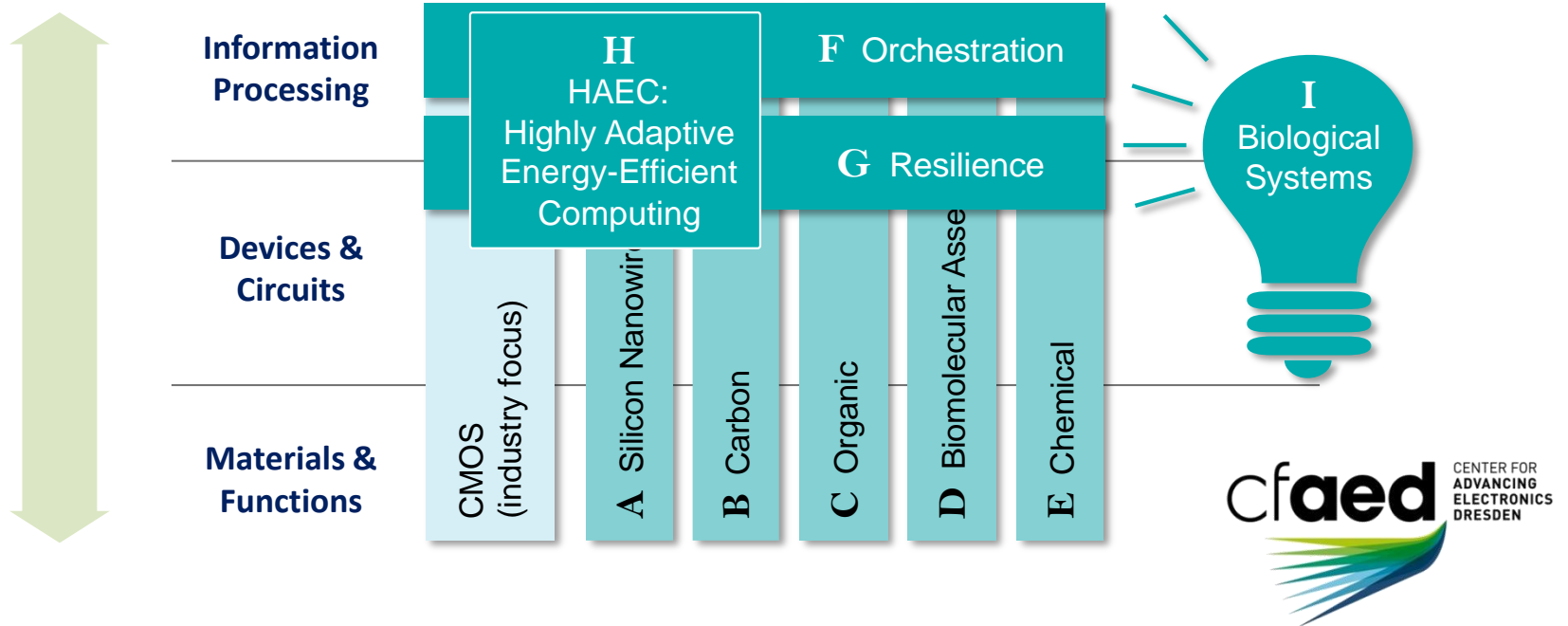
Parallel Running Query



- reduce hops by wireless „point-to-point“ interconnects
- operator → data / data → operator shipping
- adaptive hardware components (e.g., trade cache for compute units)



„More-Shots-on-Goal“





Adaptive Query Processing on Prefix Trees

Wolfgang Lehner

Fachgruppentreffen, 22.11.2012 – TU München