

# Rack-scale Data Processing System

Jana Giceva, Darko Makreshanski, Claude Barthels, Alessandro Dovis, Gustavo Alonso

Systems Group, Department of Computer Science, ETH Zurich



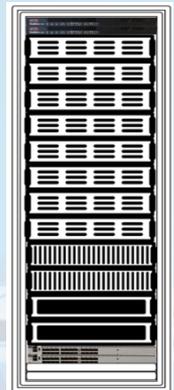
# Rack-scale Data Processing System

Jana Giceva, Darko Makreshanski, Claude Barthels, Alessandro DAVIS, Gustavo Alonso

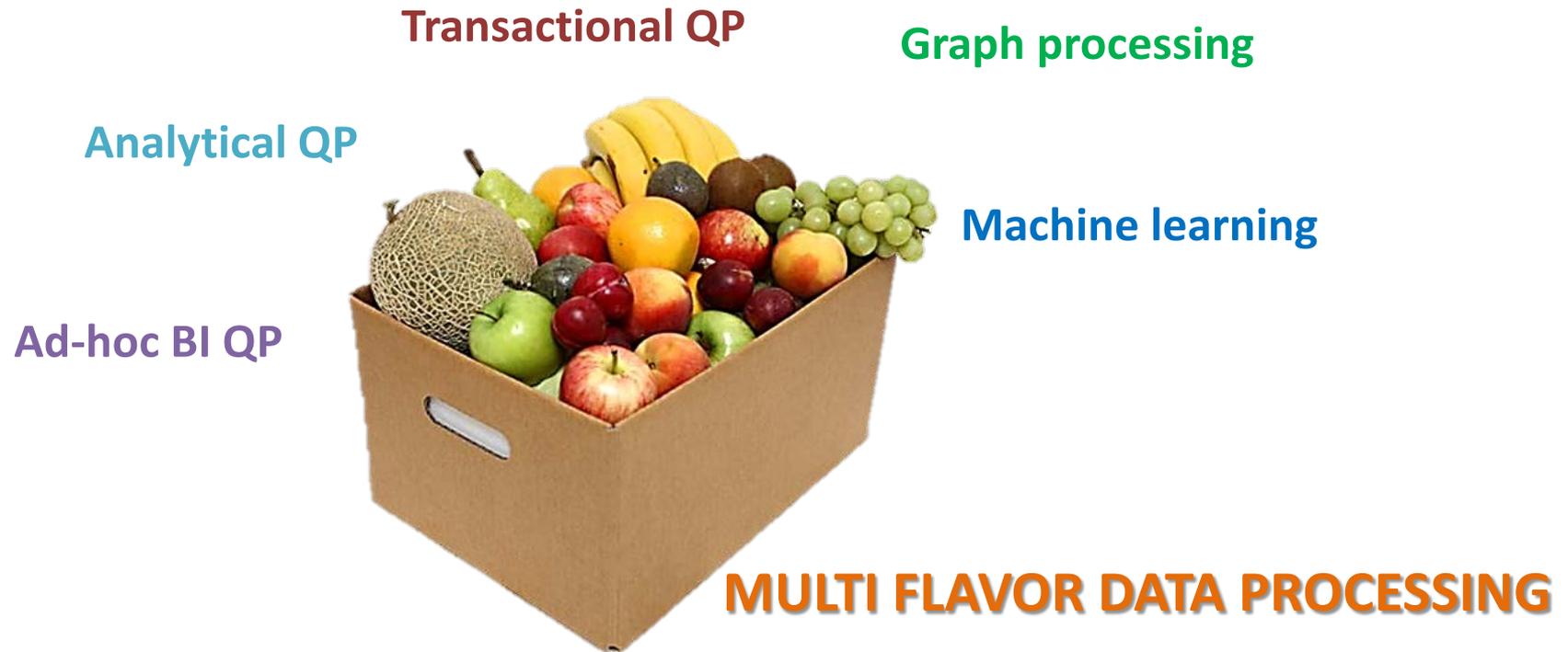
Systems Group, Department of Computer Science, ETH Zurich



## Application's perspective of a rack



# FruitBox – a data processing system



# FruitBox

- Building a system for multi-flavor data processing:
  1. Hardware that meets the resource demand.
  2. System architecture to support workload heterogeneity.
  3. Aim for 10s-100s millions of requests per second.
  4. Efficient resource utilization.

## MULTI FLAVOR DATA PROCESSING



Graph processing

Machine learning

Transactional QP

Analytical QP

Ad-hoc BI QP

# FruitBox – a **rack-scale** data processing system

- Which box could run such a heterogeneous WL?
  - A multicore is not enough
- A rack-scale system:
  - More resources
  - Better isolation
  - Blurring the machine-cluster boundaries

## MULTI FLAVOR DATA PROCESSING



Graph processing

Machine learning

Transactional QP

Analytical QP

Ad-hoc BI QP

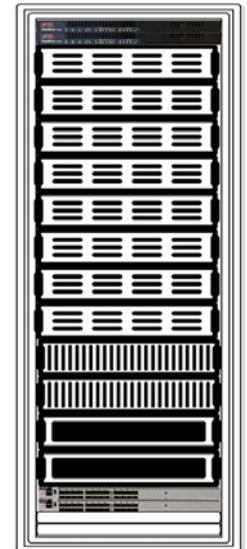


## RACK-SCALE SYSTEM

1000s of cores

TBs of RAM

InfiniBand



# Rack-scale data processing system



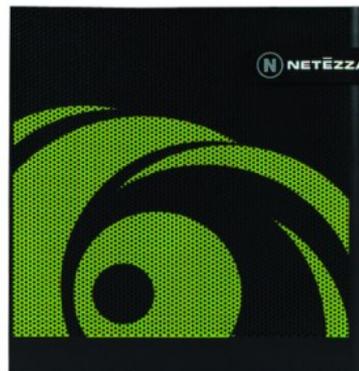
Custom build a rack-scale system  
for data processing?

Many such commercial systems exists – **Data Appliances**

**ORACLE® Exadata**



Netezza (**IBM**) TwinFin



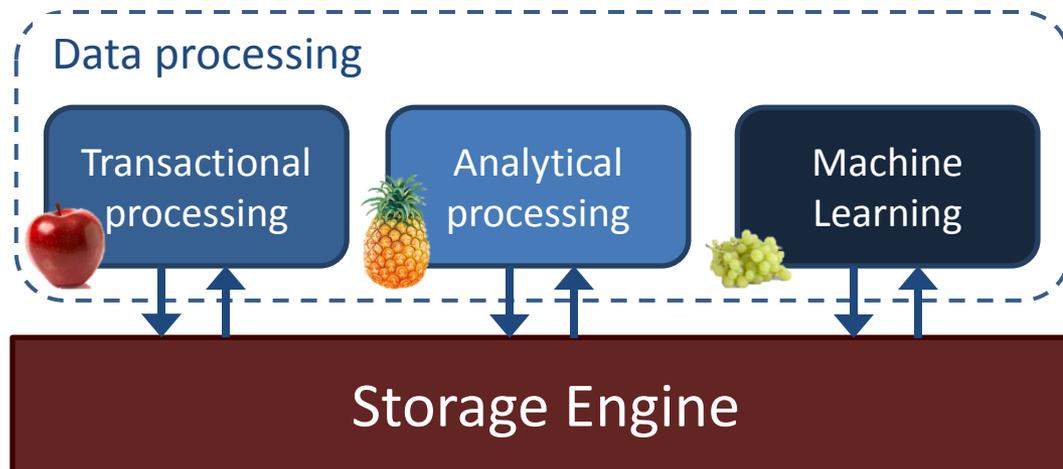
**HUAWEI** **SAP HANA**



and many more ...

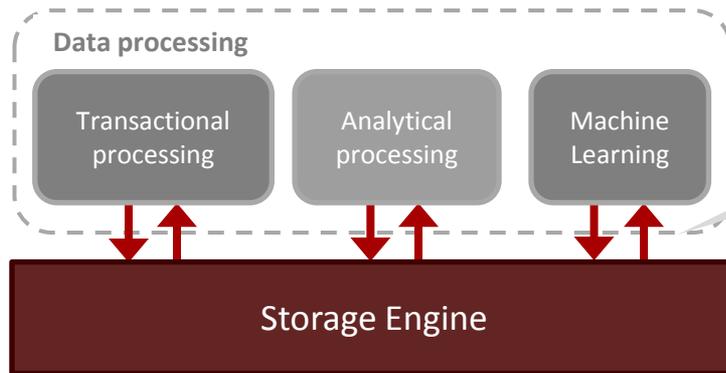
# System design for Multi-flavor data processing

- Separate data-storage from data-processing



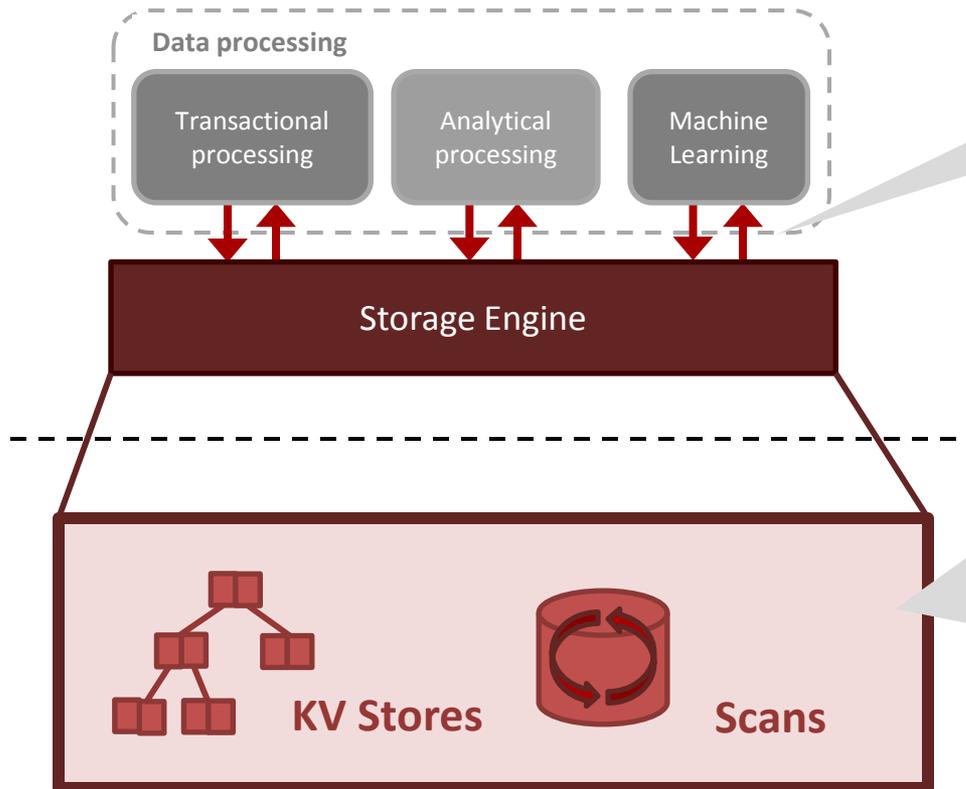
- Achieve both *physical* and *logical* data independence

# Storage Engine



*Tuple-* and *batch-based* interface to the storage engine.

# Storage Engine



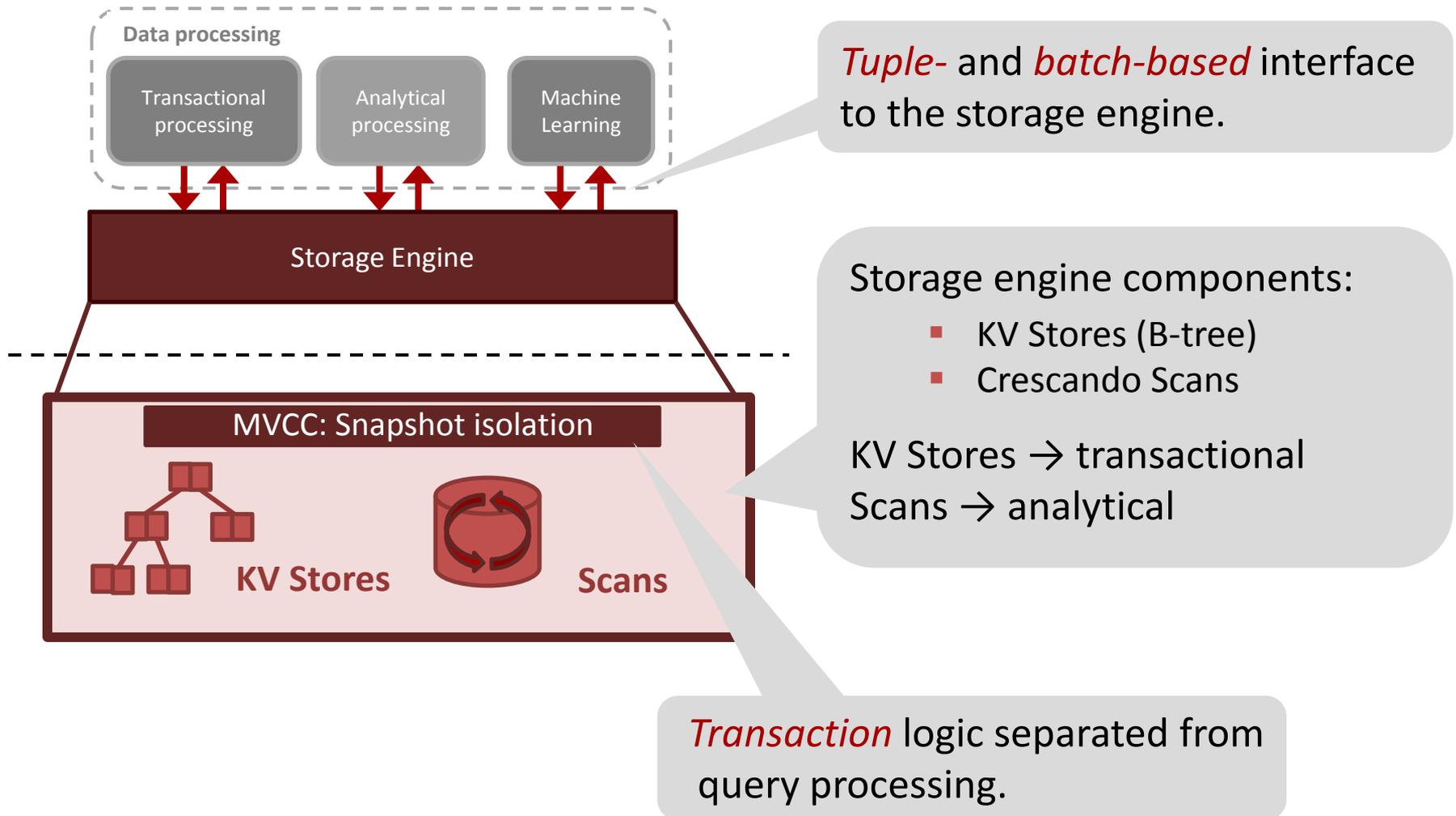
*Tuple-* and *batch-based* interface to the storage engine.

Storage engine components:

- KV Stores (B-tree)
- Crescendo Scans

KV Stores → transactional  
Scans → analytical

# Storage Engine



# Handling millions of requests/second

- It makes no sense to process them individually if they access the same data.
- Why should each query scan a TB of data?



VS



- *Batch requests* – share data, computation, bandwidth ... for higher throughput and predictable performance trading off a bit of latency.

# Efficient resource utilization

## MULTI FLAVOR DATA PROCESSING



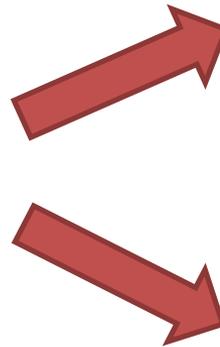
Graph processing

Machine learning

Transactional QP

Analytical QP

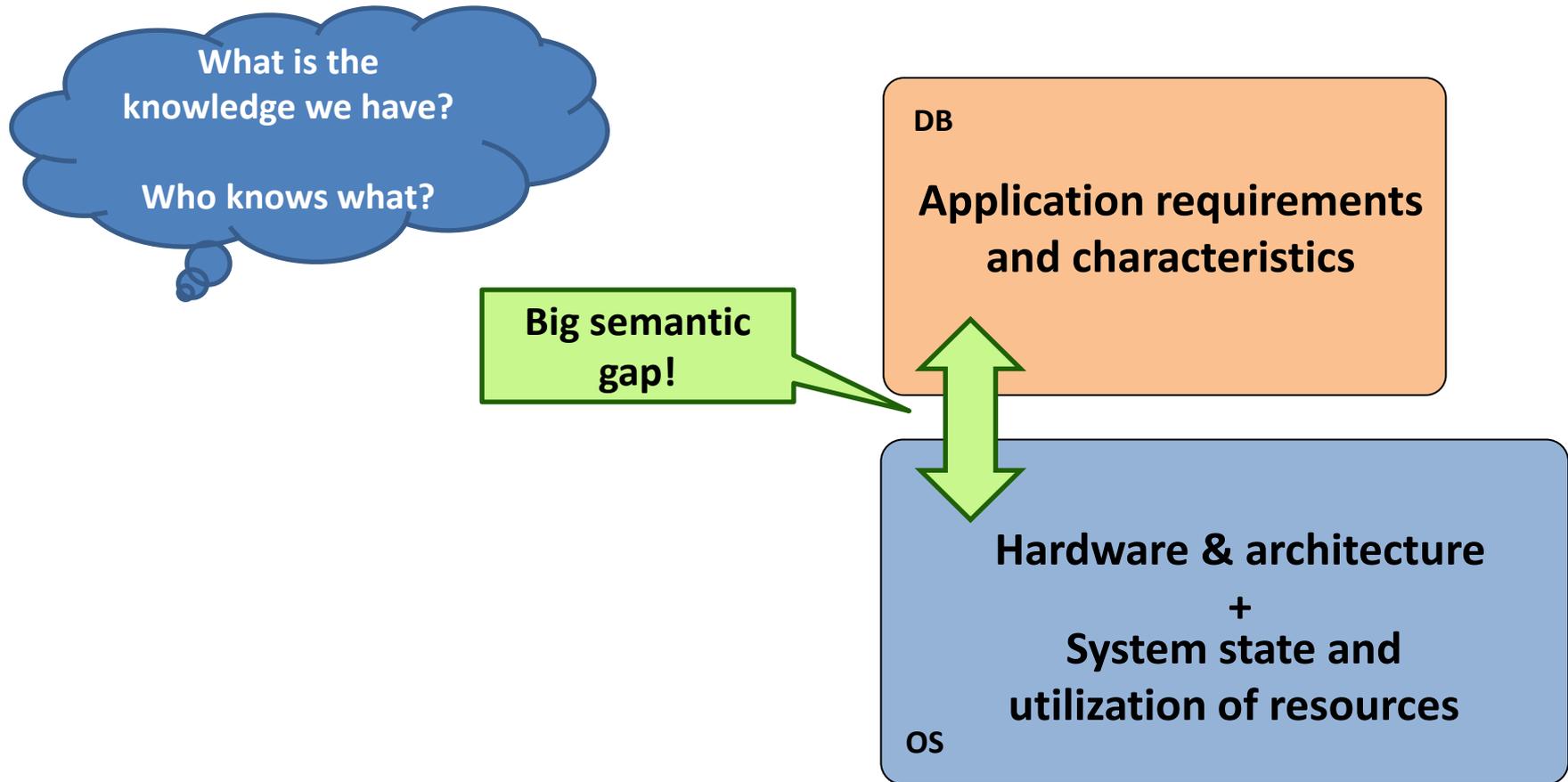
Ad-hoc BI QP



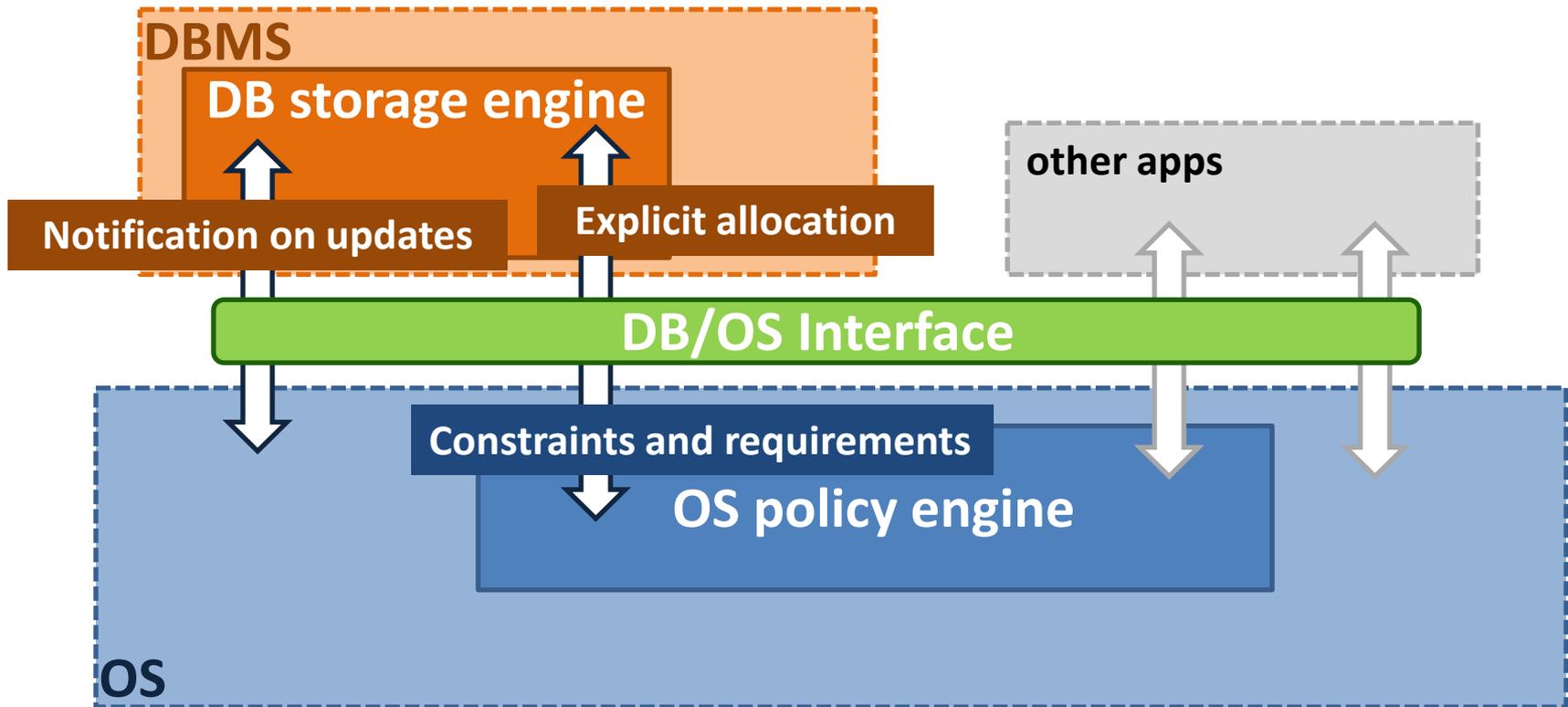
- Noisy system environment
- Load interaction
- ☹ Unpredictable performance
- ☹ Not meeting SLAs
- Resource overprovisioning
- ☹ Inefficiency and higher cost

- Getting the most out of such a complex system requires cross-layer optimization.
  - e.g. DB/OS co-design
- Already some work on multicore systems.

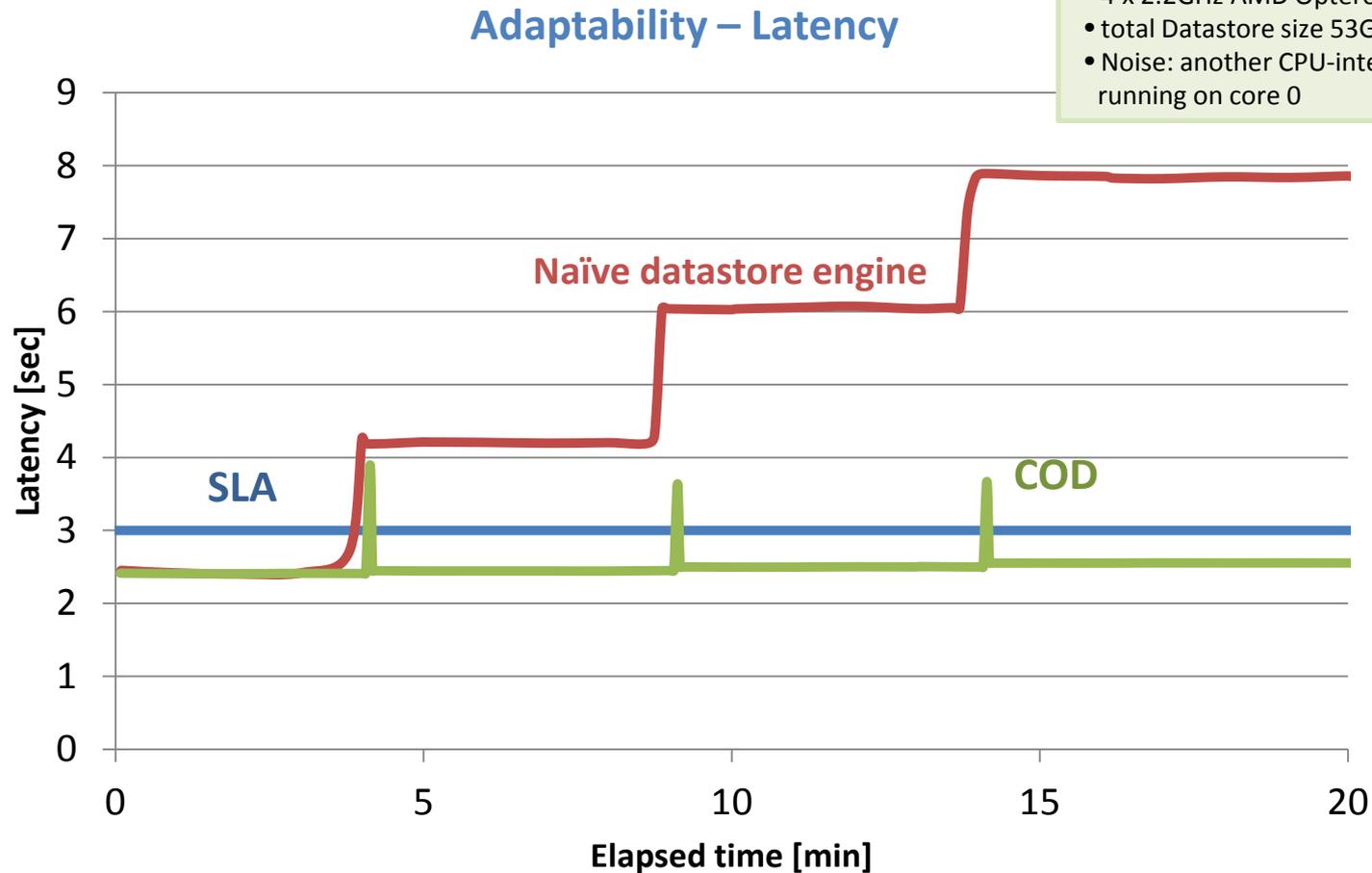
# COD: DB/OS co-design



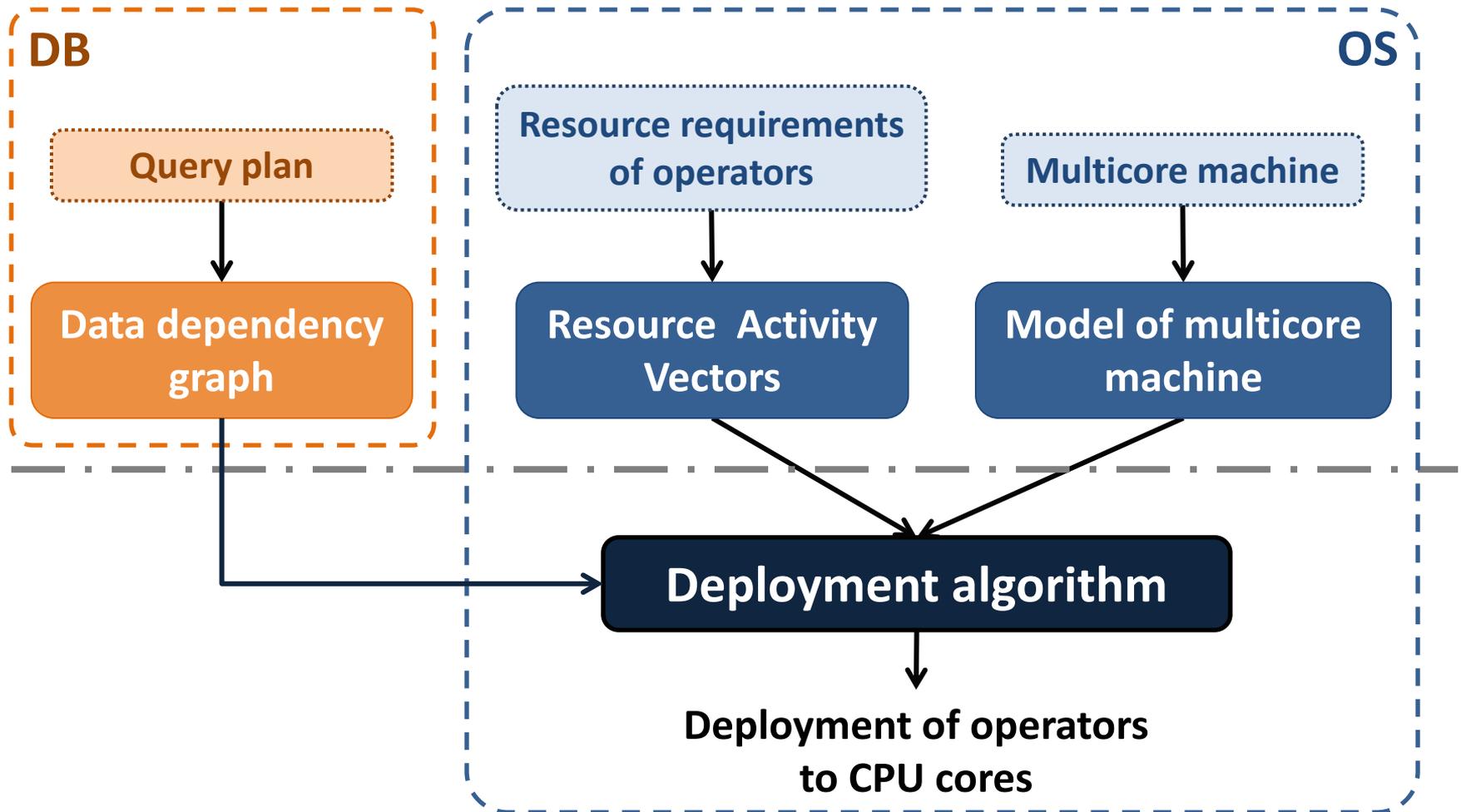
# COD's interface



# Adaptability to dynamic system state



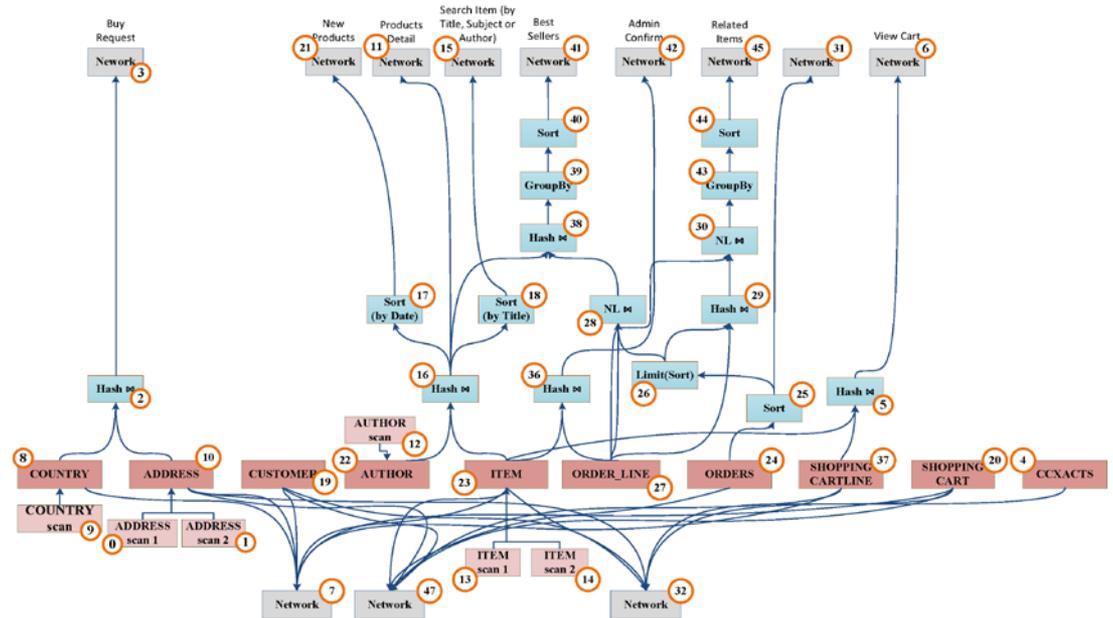
# Resource efficient deployment



# Evaluation

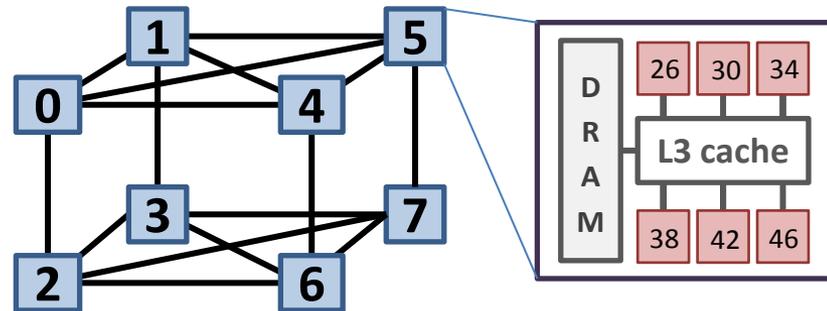
## Query plan

- SharedDB's TPC-W [1]
- 11 web-interactions in one query plan
- 44 operators
- 20GB dataset



## AMD Magnycours

- 4 x 2 dies:
  - 6 cores
  - 5 MB L3 cache
  - 16 GB NUMA node



# Comparison with standard approaches

Approaches	# cores	Throughput [WIPS]		Response Time [ms]		
		Average	Stdev	50 <sup>th</sup>	90 <sup>th</sup>	99 <sup>th</sup>
Default OS	48					
Operator per core	44					
Deployment algorithm						

# Comparison with standard approaches

Approaches	# cores	Throughput [WIPS]		Response Time [ms]		
		Average	Stdev	50 <sup>th</sup>	90 <sup>th</sup>	99 <sup>th</sup>
<b>Default OS</b>	48	317.30	31.11	8.22	72.43	82.03
<b>Operator per core</b>	44	425.86	54.34	14.59	22.93	36.08
<b>Deployment algorithm</b>						

# Comparison with standard approaches

Approaches	# cores	Throughput [WIPS]		Response Time [ms]		
		Average	Stdev	50 <sup>th</sup>	90 <sup>th</sup>	99 <sup>th</sup>
<b>Default OS</b>	48	317.30	31.11	8.22	72.43	82.03
<b>Operator per core</b>	44	425.86	54.34	14.59	22.93	36.08
<b>Deployment algorithm</b>	6					

# Comparison with standard approaches

Approaches	# cores	Throughput [WIPS]		Response Time [ms]		
		Average	Stdev	50 <sup>th</sup>	90 <sup>th</sup>	99 <sup>th</sup>
<b>Default OS</b>	48	317.30	31.11	8.22	72.43	82.03
<b>Operator per core</b>	44	425.86	54.34	14.59	22.93	36.08
<b>Deployment algorithm</b>	6	428.07	32.80	15.36	23.73	36.13

# Comparison with standard approaches

Approaches	# cores	Throughput [WIPS]		Response Time [ms]		
		Average	Stdev	50 <sup>th</sup>	90 <sup>th</sup>	99 <sup>th</sup>
Default OS	48	317.30	31.11	8.22	72.43	82.03
Operator per core	44	425.86	54.34	14.59	22.93	36.08
Deployment algorithm	6	428.07	32.80	15.36	23.73	36.13

➔ *Performance / Resource efficiency savings of x7.37*

# Conclusion



## Multi-flavor data processing system

- We have all the pieces of the puzzle



Putting them together opens a lot of opportunities.

# Conclusion



## Multi-flavor data processing system

- We have all the pieces of the puzzle

- Intelligent storage engine:
  - Co-processors, active-memory, hardware specialization (FPGAs)
- Optimizing the network stack:
  - ... for different memory access patterns
- Extend the cross-layer interface:
  - DB optimizer that is aware of the complexity of the rack
- Rack-scale resource management

Putting them together opens a lot of opportunities.