# The Case for Heterogeneous HTAP

Raja Appuswamy, Manos Karpathiotakis, Danica Porobic, and Anastasia Ailamaki Data-Intensive Applications and Systems Lab FPFI



#### HTAP – the contract with the hardware

Hybrid OLTP & OLAP Processing

High-throughput OLTP

Low-latency OLAP

#### HTAP on multicores

Massive parallelism => high concurrency Global shared memory => data sharing System-wide coherence => synchronization





# Shifting hardware landscape (1): Specialization of CPUs

#### Multisocket multicores

Intel SCC, ARM v8, Cell SPE



#### CPUs: general-purpose → customizable features



Pascal

# Shifting hardware landscape (2): Generalization of GPUs



#### GPUs: Niche accelerators $\rightarrow$ general-purpose processors



# Emerging hardware: Revisiting the contract

Current Emerging hardware	HTAP software
<ul> <li>Homogeneous Heterogeneous parallelism</li> <li>Task-parallel CPUs</li> <li>Data-parallel GPUs</li> </ul>	<ul> <li>Cannot exploit heterogeneity</li> <li>HTAP across processors</li> </ul>
<ul> <li>System-wide Relaxed cache coherence</li> <li>OS (FOS), FS (Hare)</li> <li>runtimes (Cosh)</li> </ul>	<ul> <li>Shared-everything OLTP: N/A</li> <li>No synch. sans coherence</li> </ul>
<ul> <li>Global shared memory</li> <li>Unified address space</li> </ul>	<ul> <li>Server as distributed system</li> <li>Fails to exploit shared memory</li> </ul>
Clean slate redesign in order	



## Heterogeneous HTAP (H<sup>2</sup>TAP): Caldera

- Store data in shared memory
- Run OLTP workloads on task-parallel archipelago
- Run OLAP workloads on data-parallel archipelago



Loose job-to-core assignment exploits heterogeneity



# H<sup>2</sup>TAP Challenges

- Store data in shared memory
  - Choose optimal data layout
- OLTP on task-parallel archipelago
  - Make up for (lack of) cache coherence
- OLAP on data-parallel archipelago
  - Share transactionally-consistent snapshots across processors





## Data layout

- Need to minimize PCIe data transfer to GPU
- Data access on GPU should be sequential to enable "coalescing"
- Caldera implements NSM, DSM, and PAX



PAX fits GPUs best (PCIe & coalesced accesses)



## OLTP without cache coherence

- Use Data-Oriented Transaction Execution principles
  - Thread-to-data assignment leads to partitioned data, metadata (2PL, index)





## OLTP without cache coherence

- Use explicit messaging instead of implicit latching
- Exploit shared memory by exchanging pointers instead of data





## Transactionally-consistent data sharing

- Data sharing across workloads
  - Use Unified Virtual Addressing (UVA) for CPU—GPU sharing
- Consistent data sharing via hardware snapshotting (ex: Hyper)
  - CUDA runtime restricts use in H<sup>2</sup>TAP context
- Caldera supports lightweight software snapshotting
  - OLAP queries run on immutable snapshot
  - Copy-on-write performed by update transactions

Snapshots across GPU-CPU archipelagos







#### Experiments

Setup

- Two 12-core Intel Xeon E5-2650L v3 CPUs, 256GB RAM
- GeForce GTX 980 GPU (PCIe 3.0) with 4GB memory
- TPC-C, TPC-H, YCSB in various scale factors
- Silo, MonetDB, DBMS-C

Goals

- Message passing and Software snapshotting overhead
- PAX performance compared to NSM and DSM on GPUs
- Caldera performance compared to state-of-the-art



## OLTP throughput





#### OLAP response time (incl. data movement)





#### Impact of snapshotting





#### Impact of data layout

1 table (*i1 integer, i2 integer, .... i16 integer*) SELECT SUM(colA + colB) FROM table





#### Conclusion

- Hardware architecture is changing
  - New opportunities: massive parallelism, fast interconnects
  - New challenges: heterogeneity, relaxed coherence
- Databases can and should exploit hardware trends
  - Exploit hardware heterogeneity in their core architecture design
  - Decouple system-wide coherence from shared memory
- Time to move from HTAP to H<sup>2</sup>TAP
  - H<sup>2</sup>TAP architecture: revisit age-old h/w—s/w contract
  - Caldera: Preliminary prototype to prove that H<sup>2</sup>TAP is possible