DAM◯N
*Data Management*
*On New Hardware*

# Panel discussion

## Peter Boncz (CWI)

*Architecture-Conscious Databases:*

*sub-optimization or the next big leap?*

# Sub-Optimization or Big Leap?



- The same stuff all over again..

s/Buffer Cache/L2/g

s/Page Fault/Cache Miss/g

s/Disk Block/Cache Line/g

- The final benefits are only a few percentage points
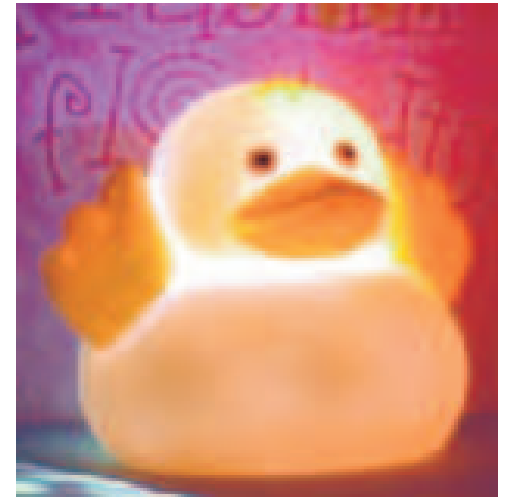
# Sub-Optimization or Big Leap?

■ Optimizing { cache use, IPC, .. } use can make a huge difference

■ Future even more interesting:

e.g. after hitting the memory wall, we now hit a CPU frequency wall.

"Computer architectures will fundamentally change""

➔ need strong help from software (e.g. multicore)

# Questions

- what will computer architectures look like in 5 years?

- do computer architecture trends/changes force us to re-think the classical DBMS architecture?

- to what extent are CPU manufacturers willing to listen to DBMS people?

- what architecture-conscious HYPEWORD data-management challenges/opportunities do you see in the next 5 years?

  HYPEWORD in { XML, stream, mobile/ubiquitous, sensor, data mining, multimedia, biological }

# The panelists



- Doug Carmean (Intel)
- Bradley Kuszmaul (MIT)
- Jignesh Patel (Univ. Michigan)
- Babak Falsafi (CMU)
- Kenneth Ross (Columbia Univ.)

# The panelists



- Doug Carmean (Intel)
- **Bradley Kuszmaul (MIT)**
- Jignesh Patel (Univ. Michigan)
- Babak Falsafi (CMU)
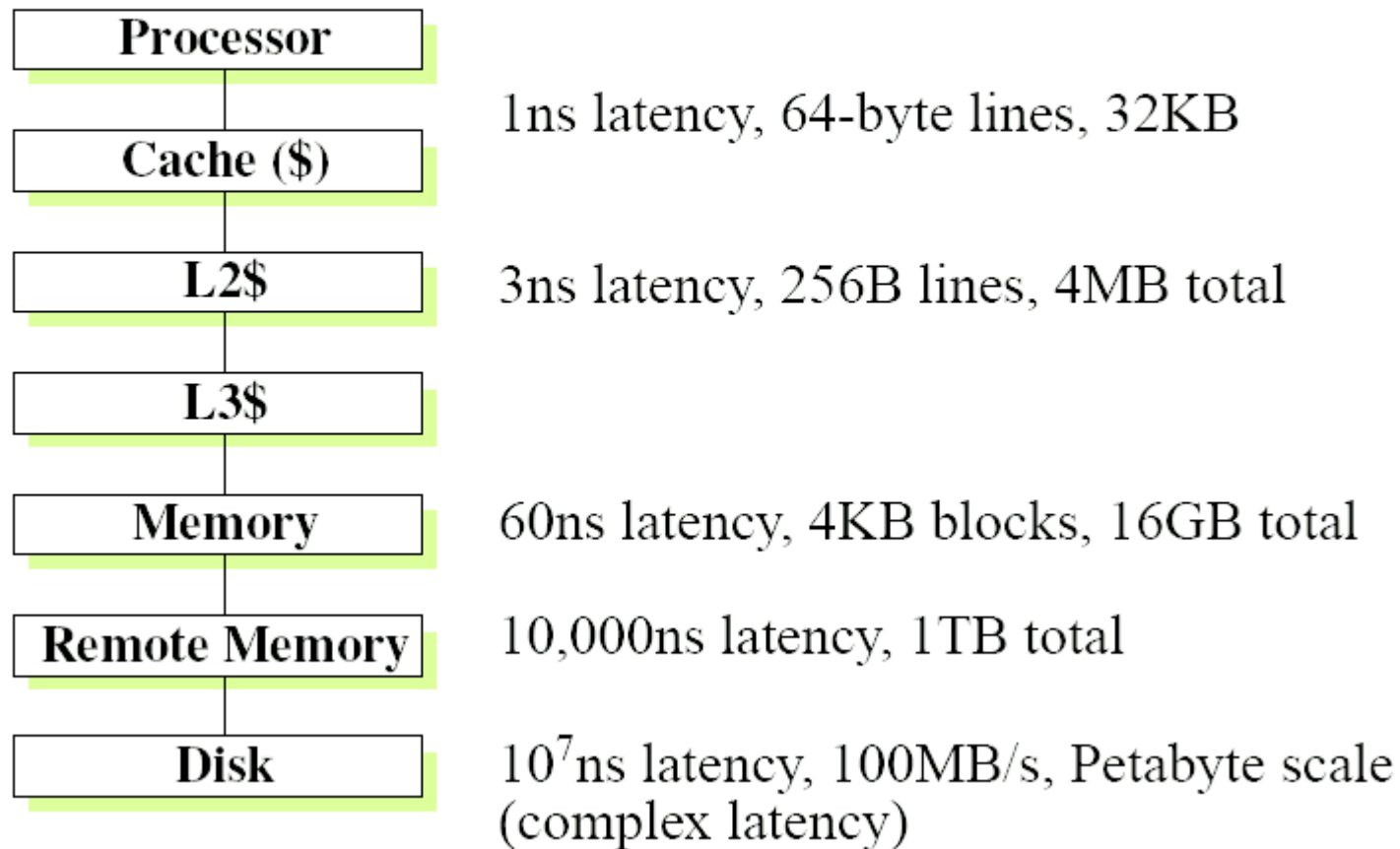- Kenneth Ross (Columbia Univ.)

# Cache-Oblivious Data Structures for Databases
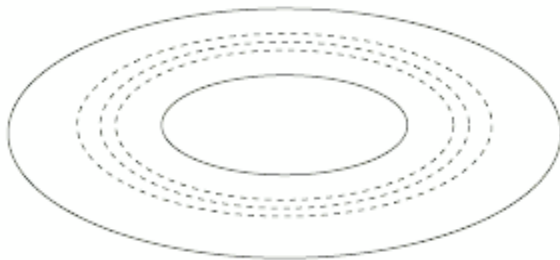## Bradley C. Kuszmaul
## MIT CSAIL

The cache hierarchy is where the performance goes.

# The Memory Hierarchy

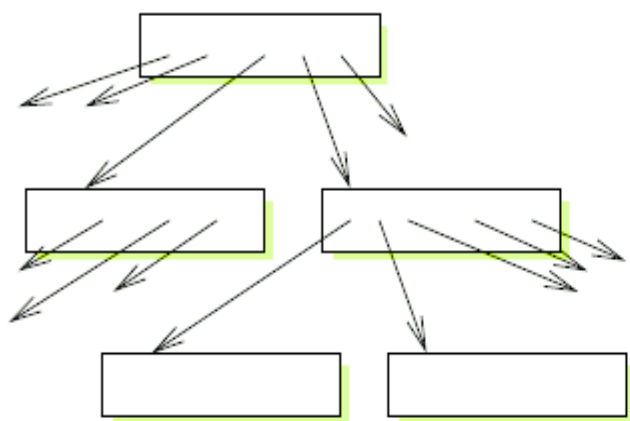| | |
|---|---|
| **Processor** | |
| **Cache ($)** | 1ns latency, 64-byte lines, 32KB |
| **L2$** | 3ns latency, 256B lines, 4MB total |
| **L3$** | |
| **Memory** | 60ns latency, 4KB blocks, 16GB total |
| **Remote Memory** | 10,000ns latency, 1TB total |
| **Disk** | $10^7$ns latency, 100MB/s, Petabyte scale (complex latency) |

# Disk Systems Are Complex



- Rotational latency 5ms,
- long seek 8ms,
- short seek 1ms.
- Outer tracks contain twice as much data as inner tracks, with the same transfer time (1 rotation).
- Caching in the disk controller.

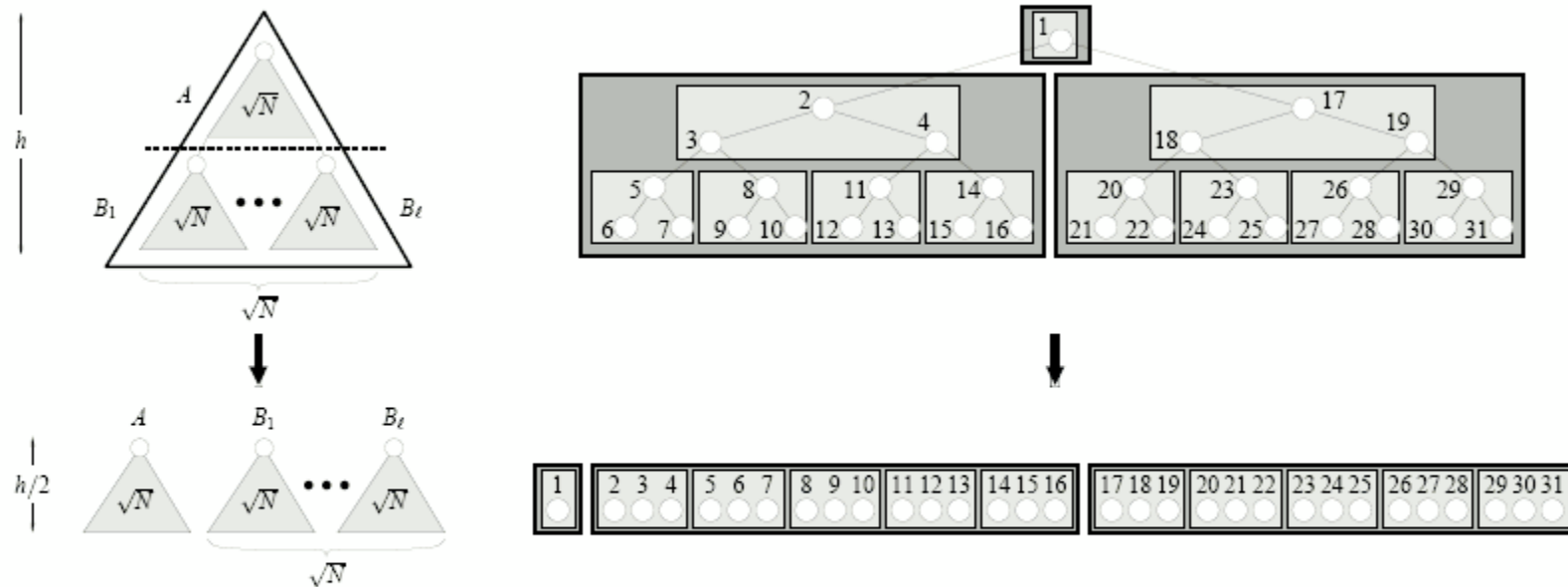Half-power point is size at which half your cost is latency,

# B-Tree Block Size

- B-trees choose the node degree so that each node fits in a block.
- What block size to use?
  - 64B (cache lines),
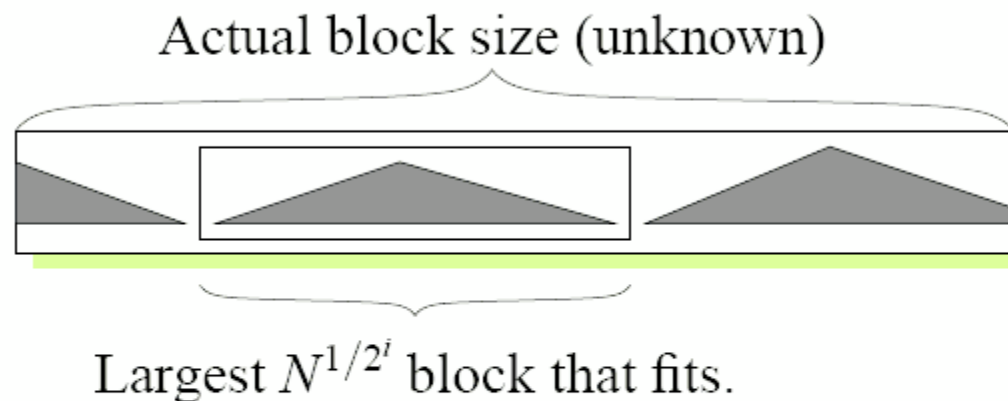  - 4KB (page size),
  - $1/3$MB (random seek half-power point)

Small blocks make poor use of disk. Big blocks make poor use of cache (and may be expensive for insertions.)

# A Static Cache-Oblivious B-tree



First lay out the tree into blocks of size $\sqrt{N}$. Then lay out those blocks into blocks of size $N^{1/4}$. (And those into blocks of size $N^{1/8}$, ...)

# Cache-Oblivious B-trees Optimize for All Block Sizes Simultaneously

Actual block size (unknown)



Largest $N^{1/2^i}$ block that fits.

Two blocks of size $\sqrt{B}$ give the same fanout as one block of size $B$. $\rightarrow$ At most 2X more block transfers than optimal.

(Potentially another 2X due to alignment. But small blocks, which pay the first factor of two, are unlikely to cross the actual block boundary.)

This analysis applies to *all* block sizes at the same time.

# Some Performance measurements

Random searches:

| Data structure | Average time per search | |
|---|---|---|
| | small-machine | big-machine |
| CO Btree | 12.3ms | 13.8ms |
| Btree: 4KB Blocks: | 17.2ms | 22.4ms |
| 16KB blocks: | 13.9ms | 22.1ms |
| 32KB blocks: | 11.9ms | 17.4ms |
| 64KB blocks: | 12.9ms | 17.6ms |
| 128KB blocks: | 13.2ms | 16.5ms |
| 256KB blocks: | 18.5ms | 14.4ms |
| 512KB blocks: | | 16.7ms |

Range queries run 100 times faster.

Conclusion: Must optimize for the cache heirarchy. Can do so in an architecturally independent way.

# The panelists



- Doug Carmean (Intel)
- Bradley Kuszmaul (MIT)
- **Jignesh Patel (Univ. Michigan)**
- Babak Falsafi (CMU)
- Kenneth Ross (Columbia Univ.)

# Architecture-Conscious Databases: Sub-optimization or the Next Big Leap?

## Jignesh M. Patel
## University of Michigan

**Those who cannot remember the past are condemned to repeat it.**
**- George Santayana**

# History (from a DB perspective)

- 1970s, 80s: Birth of RDBMS
  - Disk I/O are expensive: buffer management, prefer sequential IOs.
  - Customized hardware for DBMS?
    - Customized hardware is not cost-effective or portable, better to use smart software-based solutions for query processing.
- 1980s, 90s: High performance and scalable DBMS
  - Parallelism: RDBMS are very amenable to parallelism.
    - Shared-nothing - essentially application level-parallelism.
    - Need to worry about startup, interference, and skew.
  - Expanding DBMS functionality: Object*, external functions, rules, …
- 1990s, 21st century: Ubiquitous DBMS
  - Expanding DBMS applications: streams, scientific, semi-structured, personal data management, …
  - Architecture-conscious RDBMS:
    - Memory Wall: Design cache-aware methods, prefetching, …
    - Exploit new processor features: SIMD, co-processors (GPU), …

# Future Computer Architecture and impact on DBMS

- ## SMP on a chip

  - ► Shared-nothing parallel DBMS will be the default installation.

- ## Memory hierarchy will continue

  - ► But processor clock speeds are not increasing rapidly

# CPU extensions for DBMSs

- No compelling reason to add DBMS-specific extensions

- Recall the database machine's era

  ▶ DBMS has smart and efficient software-based techniques

- What is the payoff for the hardware vendor?

  ▶ Servers: Application servers (Java + DB), file servers, ...

  ▶ Clients: Entertainment, Virus/spyware scanners, personal data management, ...

- DBMS part of a complex suite of software even on server machines

  ▶ Lots of time is spent in external function and (Java) application code.

# Rethink traditional DBMS Architecture?

- We already have shared-nothing parallelism.
- Memory hierarchy fundamentally remains the same: smaller and faster memory is closer to the processor
  - ► For shared-nothing parallelism, synchronization may be cheaper with shared on-chip memory
- Coprocessors for RDBMS operations

**NO!**

  - ► There will be less pressure to use coprocessors.
  - ► Not clear if this approach has a significant performance advantage when compared to the **best** current methods on regular processors.
  - ► End-to-end application performance?
    - ► RDBMS query processing will be a shrinking component.
  - ► Speed is only a small factor. Portability is important.

# New Challenges

## The role of DBMS is expanding!

- Efficient fuzzy matching algorithms on text, personal digital records, multimedia, graphs.

- Data mining will be more common as data volume continues to explode and complexity of analysis increases. Relatively little work on parallel algorithms here.

- Scientific workloads: Lots of very computationally expensive data analysis. Needs massively parallel methods. Data volume is exploding.

- But need clean programming interfaces.

# The panelists



- Doug Carmean (Intel)
- Bradley Kuszmaul (MIT)
- Jignesh Patel (Univ. Michigan)
- **Babak Falsafi (CMU)**
- Kenneth Ross (Columbia Univ.)

DaMoN 2005 panel on

# Architecture-Conscious Databases: sub-optimization or the next big leap?

*Babak Falsafi*

**Computer Architecture Lab**
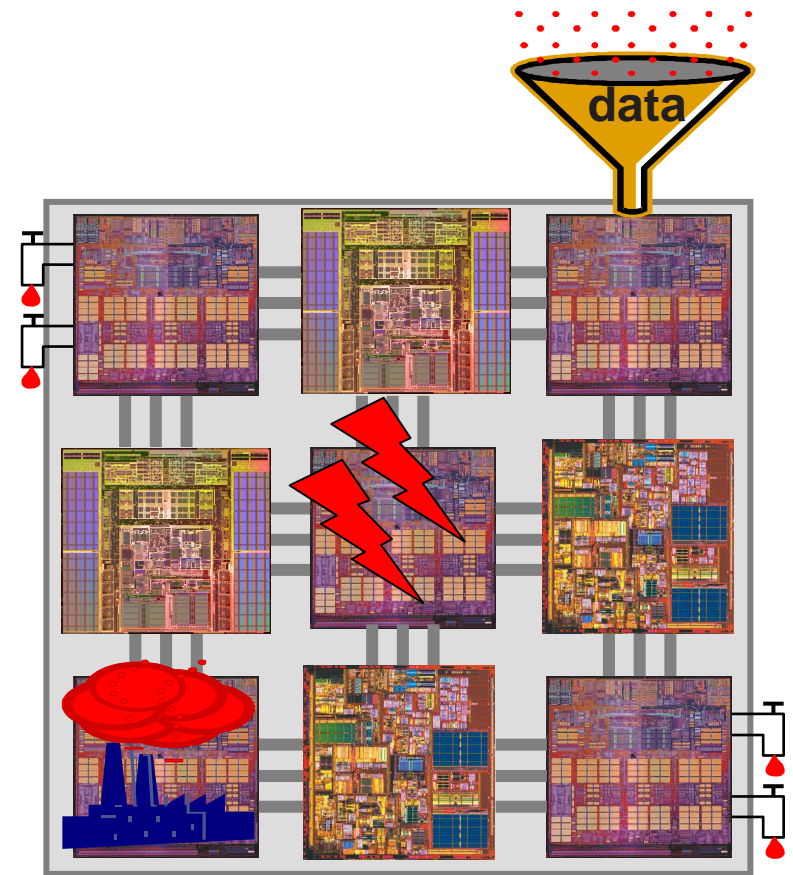**Carnegie Mellon**
http://www.ece.cmu.edu/CALCM

# Technology Scaling Trends

Good news:
- 100B trans/chip by 2015
- Tens of cores

Bad news:
- Heterogeneous, smoking and unreliable cores
- Faster but constrained clocks
- Off-chip latency/bw bottlenecks continue



**2015 Multi-core Chip**

# Scaling Implications for DB Servers

## Lots of cores & on-chip memory:

– Parallelism galore: intra-transaction

– Sea of memory, not hierarchies

## Heterogeneous resources:

– ILP, TLP, and DLP (vector)?

– Both HW-/SW-managed caches

## Unreliable HW:

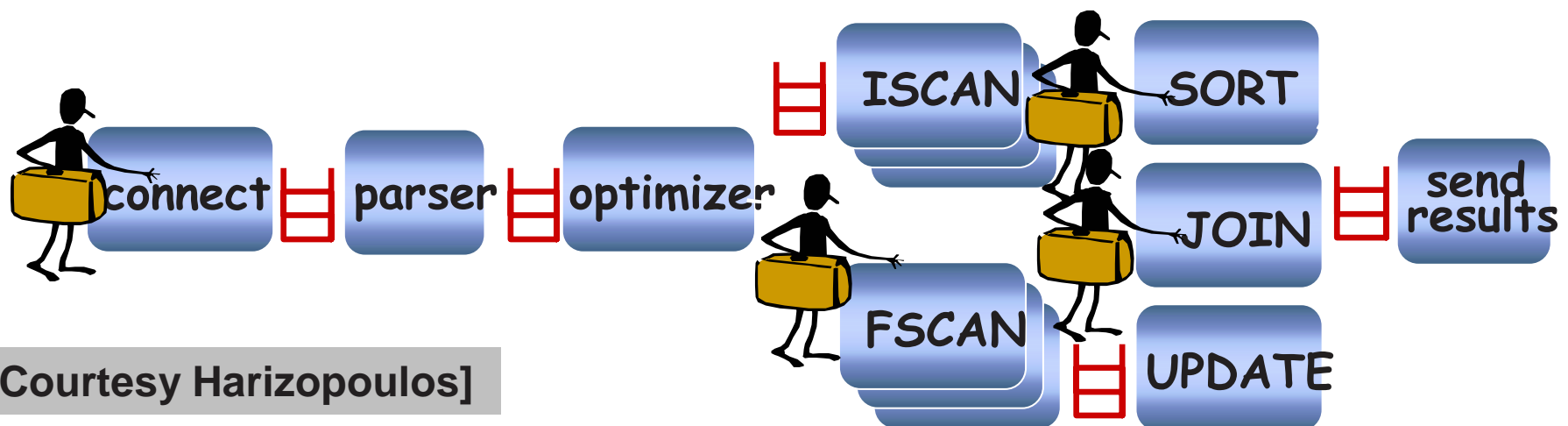– Fine-grain transaction semantics

**Need tight HW/SW collaboration!**

# One Solution:
# Staged Server Architecture

Componentize server into operators

Operator-level pipelining

– Map operators to cores

– HW/SW data streaming
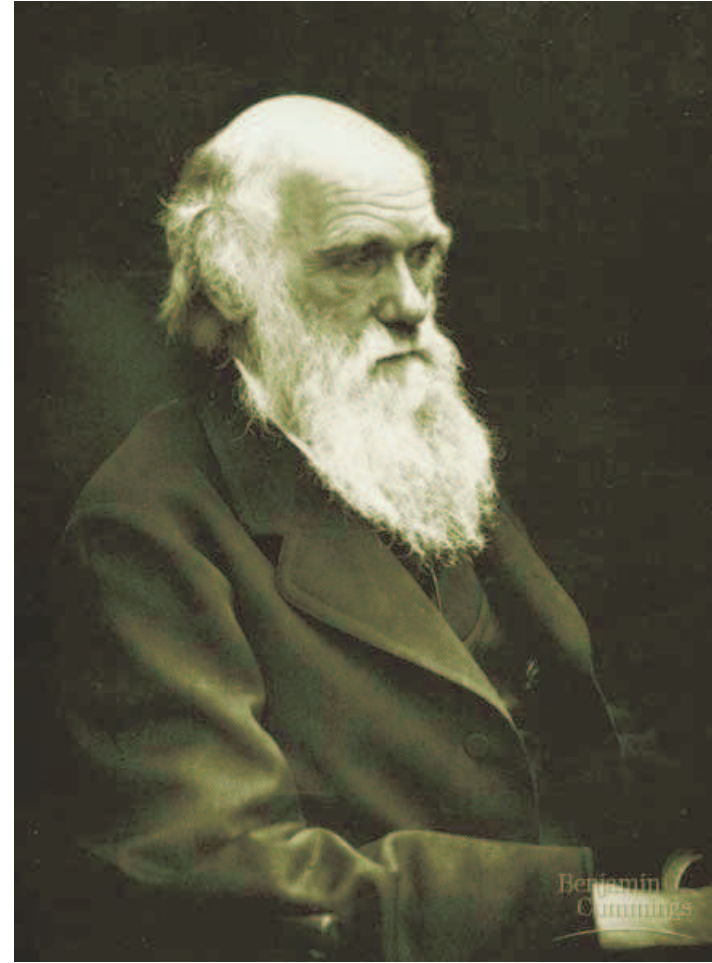


[Courtesy Harizopoulos]

# The panelists



- Doug Carmean (Intel)
- Bradley Kuszmaul (MIT)
- Jignesh Patel (Univ. Michigan)
- Babak Falsafi (CMU)
- **Kenneth Ross (Columbia Univ.)**

# Architecture-Conscious Databases Panel

# DaMoN 2005

# Ken Ross

# Charles Darwin

# Survival of the Fittest

- Small genetic changes make the entity fitter, and are selected by the environment $\rightarrow$ Evolution

- Small $\neq$ Unimportant

- Most changes not visible

- Changes to the environment create new optimization criteria for natural selection.

# Survival of the Fittest DB

- Small implementation changes make the DB faster, and are selected by the marketplace → Evolution
- Small ≠ Unimportant
- Most changes not visible to users or developers (and shouldn't be!)
- Changes to the technology environment create new optimization criteria for DB implementation.
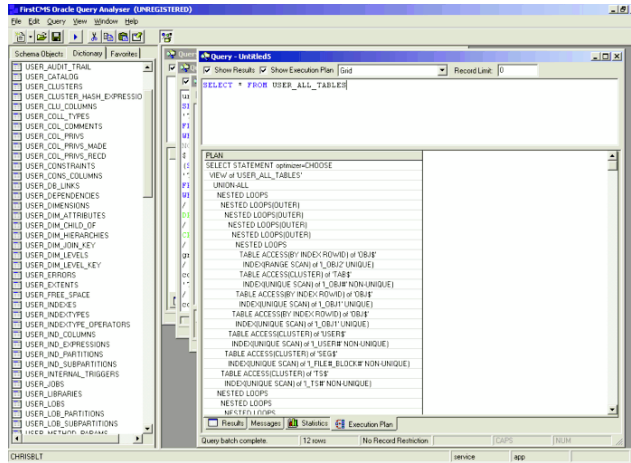
# "Sub-Optimization" or "The Next Big Thing"?

- Sub-optimization, of course!
- But many small changes can radically alter the overall structure

# The Old Days

**DB Implementor**

**DB software**



Disk block size,

Physical layout, …

Sequential access,
…

**Disk**

As disks got better, there was little need for fundamental change.

# Now

**DB software**

**CAM**

**NPU**

**DB Implementor**

!

Broadcast, …

Associative access, …

More functionality in the disk system

**Programmable Storage Subsystem**

Cache miss penalty, …

Branch mispredictions, SMT effects, …

SIMD, Parallelism, …

Streaming, geometry, …

**RAM**

**CPU**

AMD Opteron

**Compiler**

GCC

Instruction set, …

**GPU**

GEFORCE FX
nVIDIA
G30480.1 0304A2
S TAIWAN

# Challenges

- Too much is visible to the DB implementor!

- Abstractions, modularity

- Interactions: many competing "small" changes

- Management

# Summary

- [carmean] CPU future = Cell done right
- [carmean] Hardware companies are willing to listen, also on APIs
- [kuszmaul] Not cache-conscious ➔cache-oblivious (will also do disk!)
- [patel] Focus on extensibility, RDBMS query processing a shrinking component optimize/parallelize: fuzzy matching, graph algorithms, data mining, scientific
- [patel]  RDBMS architecture is fine; will handle massive multiprocessors well
- [patel] limitations of coprocessors will prevent adoption
- [falsafi] heterogeneous CPUs, sea of programmable memory (no hierarch. cache)
- [falsafi] bottomline: need for tight hardware/software collaboration
- [falsafi] solution: staged servers architecture (operator-level pipelining)
- [ross] survival of the fittest database: small changes lead evolution
- [ross] database implementor is confused (from disk alone to many arch. factors) challenge=> abstraction