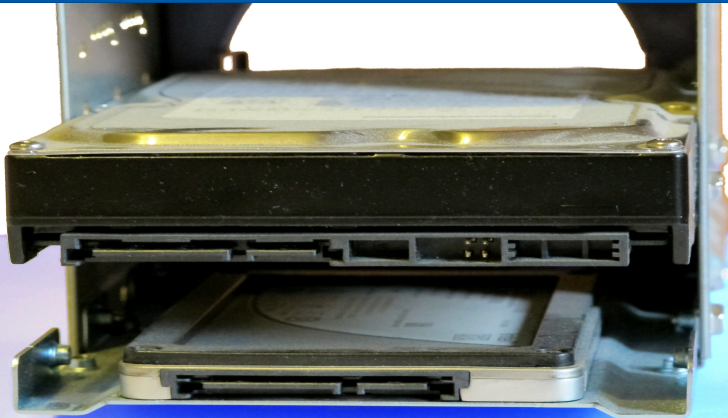# Making Cost-Based Query Optimization Asymmetry-Aware

**Daniel Bausch**, Ilia Petrov, and Alejandro Buchmann
**{bausch, petrov, buchmann}@dvs.tu-darmstadt.de**

TECHNISCHE
UNIVERSITÄT
DARMSTADT

flashyDB DVS

**Asymmetry in new storage devices**

- Writing to Flash memory is slower than reading from it
- This also applies to emerging non-volatile memories (PCM, etc.)
- Small writes to random locations on Flash are even more slow
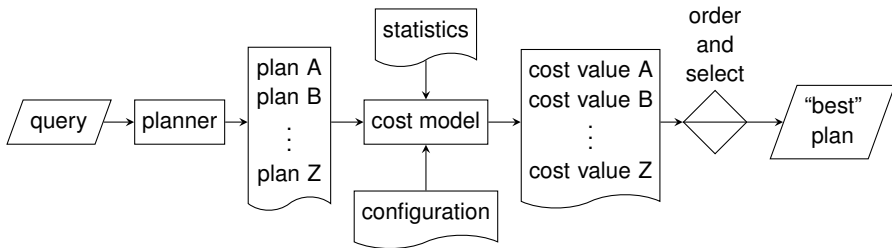- Random reads from Flash are only $\frac{1}{3}$ slower than sequential reads[1]

---

[1] on Intel X25-E using full command queue
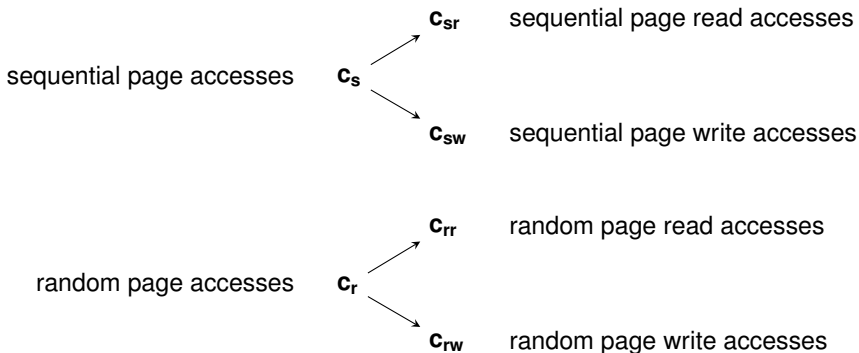
## Cost-Based Query Optimization

TECHNISCHE
UNIVERSITÄT
DARMSTADT

- estimation of run-time before real execution (e.g. in PostgreSQL)
- model comprised of functions like

$$c(seqscan) = \underbrace{\mathbf{c_s} \, \|R\|_p}_{\text{I/O cost}} + \underbrace{q_{R,0} + (\dot{\mathbf{c}}_{\mathbf{cpu}} + \dot{q}_R) \, \|R\|_t}_{\text{CPU cost}}$$

```
                          ┌───────────┐
                          │ statistics│
                          └─────┬─────┘
                                │                                    order
                                ▼                                     and
          ┌──────┐     ┌─────────────┐     ┌──────────────┐         select
┌──────┐  │plan A│     │             │     │cost value A  │        ◇         ┌───────┐
│query │→│planner│→│plan B│ → │ cost model  │ → │cost value B  │ → ◇ → │"best" │
└──────┘  │  ⋮   │     │             │     │    ⋮         │        ◇         │ plan  │
          │plan Z│     └─────────────┘     │cost value Z  │                  └───────┘
          └──────┘            ▲            └──────────────┘
                              │
                       ┌──────────────┐
                       │configuration │
                       └──────────────┘
```

flashyDB DVS

# Splitting Parameters

sequential page accesses $c_s$

$c_{sr}$  sequential page read accesses

$c_{sw}$  sequential page write accesses

random page accesses $c_r$

$c_{rr}$  random page read accesses

$c_{rw}$  random page write accesses

# Cost functions for "pure load" algorithms

TECHNISCHE
UNIVERSITÄT
DARMSTADT

| cost function | kind | original | replacement |
|---|---|---|---|
| sequential scan | read only | $c_s$ | $c_{sr}$ |
| index scan | read only | $c_s$ | $c_{sr}$ |
| | | $c_r$ | $c_{rr}$ |
| bitmap scan | read only | $c_s$ | $c_{sr}$ |
| | | $c_r$ | $c_{rr}$ |
| TID scan | read only | $c_r$ | $c_{rr}$ e |
| materialization | write only | $c_s$ | $c_{sw}$ |
| re-scan | read only | $c_s$ | $c_{sr}$ |

flashyDB ⚙ DVS

## Cost function of sort algorithm

$$c_{io}(sort) = \overbrace{2\,\|S\|_p\,\lceil log_m n\rceil\left(\tfrac{3}{4}\mathbf{c_s} + \tfrac{1}{4}\mathbf{c_r}\right)}^{\text{startup}}$$

| **blktrace stats** | write | | read | |
|---|---|---|---|---|
| | s | r | s | r |
| external sort of unordered data | | | | |
| external sort of ordered data | | | | |
| sort-merge join | | | | |

$$c_{io}(sort)_{rw} = \underbrace{\|S\|_p\left(\mathbf{c_{sw}} + \left(\lceil log_m n\rceil - 1\right)\frac{\mathbf{c_{sw}+c_{rw}}}{2} + \lceil log_m n\rceil\,\frac{\mathbf{c_{sr}+c_{rr}}}{2}\right)}_{\text{startup}}$$

# Cost function of hash join

$$c_{\mathrm{io}}(hashjoin) = \overbrace{\|P_i\|_{\mathrm{p}} \, \mathbf{c_s}}^{\text{startup}} + \left( \|P_i\|_{\mathrm{p}} + 2 \, \|P_o\|_{\mathrm{p}} \right) \mathbf{c_s}$$

| blktrace stats | write | | read | |
|---|---|---|---|---|
| | s | r | s | r |
| hash join | | | | |

$$c_{\mathrm{io}}(hashjoin)_{rw} = \underbrace{\|P_i\|_{\mathrm{p}} \, \mathbf{c_{rw}}}_{\text{startup}} + \|P_i\|_{\mathrm{p}} \, \mathbf{c_{sr}} + \|P_o\|_{\mathrm{p}} \, (\mathbf{c_{rw}} + \mathbf{c_{sr}})$$

# System and Load

- system with tight memory configuration
  $\Rightarrow$ simulate data-intensive systems under high load
- application class benchmark based on TPC-H specs

# Calibration



- based on *simulated annealing*[Kirk1983], accepts inferior configuration at decreasing probability
- modify by multiplication with logarithmic normally distributed random variable
- cooling cycle of 100 iterations, restarted 100 times
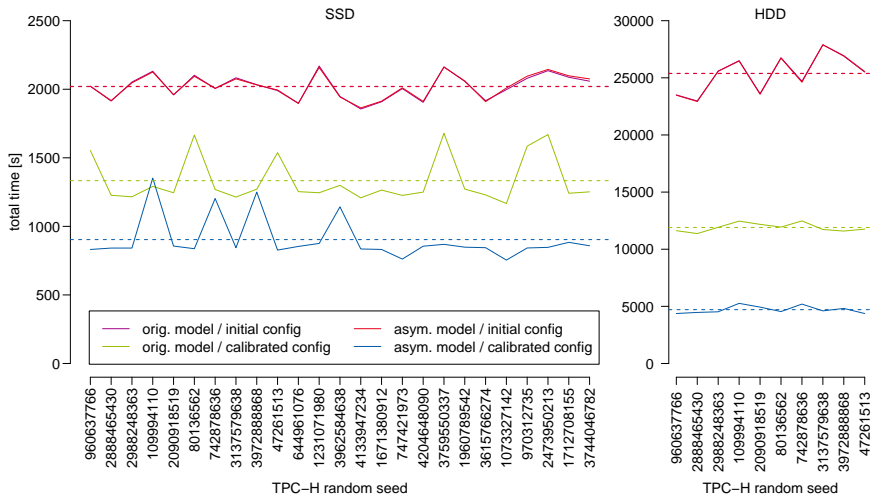
# Found "Optimal" Settings

| | SSD | | | HDD | |
|---|---|---|---|---|---|
| | original model | asymmetric model | | original model | asymmetric model |
| | $c_s$ = 1.00000 | $c_{sr}$ = 1.00000 | | $c_s$ = 1.00000 | $c_{sr}$ = 1.00000 |
| | | $c_{sw}$ = 49.91840 | | | $c_{sw}$ = 110.21139 |
| | | $c_{rr}$ = 5.62724 | | | $c_{rr}$ = 19.25494 |
| | $c_r$ = 6.77405 | $c_{rw}$ = 19.08421 | | $c_r$ = 29.04790 | $c_{rw}$ = 20.18467 |
| | $\dot{c}_{cpu}$ = 0.00121 | $\dot{c}_{cpu}$ = 0.00003 | | $\dot{c}_{cpu}$ = 0.00280 | $\dot{c}_{cpu}$ = 0.00082 |
| | $\hat{\dot{c}}_{cpu}$ = 0.03658 | $\hat{\dot{c}}_{cpu}$ = 0.01608 | | $\hat{\dot{c}}_{cpu}$ = 0.03718 | $\hat{\dot{c}}_{cpu}$ = 0.00045 |
| | $c_{op}$ = 0.00016 | $c_{op}$ = 0.00008 | | $c_{op}$ = 0.00004 | $c_{op}$ = 0.00119 |

# Comparison
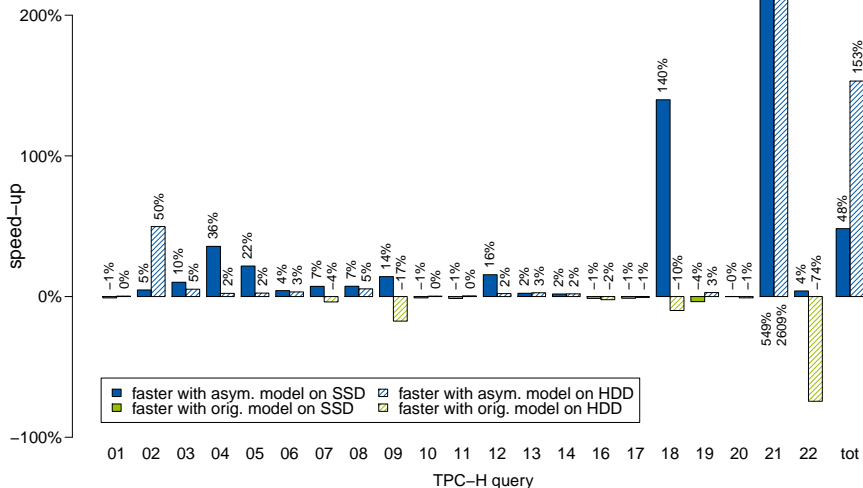
# Individual Query Speed-Up

TECHNISCHE
UNIVERSITÄT
DARMSTADT

```
                                    Nested Loop
                                    509541.375
                                   /          \
                    Nested Loop          Index Scan | using pk_orders on orders
                    102642.716                      | orders.o_orderkey=lineitem.l_orderkey
                   /         \                       0.209*1934051=404216.659
          Merge Join | partsupp.ps_partkey=lineitem.l_partkey     Index Scan | using pk_supplier on supplier
                     | Join Filter: lineitem.l_suppkey=partsupp.ps_suppkey    | supplier.s_suppkey=lineitem.l_suppkey
                       88687.896                                              0.006*1934051=11604.306
         /          \
Index Scan | using i_ps_partkey on partsupp     Materialize
             2986.769                            79808.663
```

| Nested Loop |
|---|
| 68231.426 |

| Hash Join | lineitem.l_suppkey=supplier.s_suppkey AND lineitem.l_partkey=partsupp.ps_partkey |
|---|---|
| 54756.688 | |

| Index Scan | using pk_orders on orders orders.o_orderkey=lineitem.l_orderkey |
|---|---|
| 0.008*1468202=11745.616 | |

| Seq Scan | on lineitem |
|---|---|
| 13467.469 | |

| Hash | (1 Batch) |
|---|---|
| 6460.738 | |

TECHNISCHE
UNIVERSITÄT
DARMSTADT

| Nested Loop Anti Join | Join Filter: l3.l_suppkey=l1.l_suppkey |
|---|---|
| 327228.606 | |

| Hash Semi Join | l1.l_orderkey=l2.l_orderkey Filter: l2.l_suppkey<>l1.l_suppkey |
|---|---|
| 63120.425 | |

| Index Scan | using i_l_orderkey on l3 l3.l_orderkey=l1.l_orderkey Filter: l3.l_receipdate>l3.l_commitdate |
|---|---|
| 0.362*726869=263126.578 | |

| Hash Join | l1.l_suppkey=supplier.s_suppkey |
|---|---|
| 22225.914 | |

| Hash | (128 Batches, 9271kB) |
|---|---|
| 32415.030 | |

flashyDB  DVS

# Discussion – Query 21 – the fast plan




Nested Loop Anti Join | Join Filter: l3.l_suppkey<>l1.l_suppkey
62690.508

Merge Semi Join | orders.o_orderkey=l2.l_orderkey
Join Filter: l2.l_suppkey<>l1.l_suppkey
59574.038

Index Scan | using i_l_orderkey on l3
l3.l_orderkey=l1.l_orderkey
Filter: l3.l_receiptdate>l3.l_commitdate
0.008*350514=2804.112

Merge Join | orders.o_orderkey=l1.l_orderkey
31968.016

Index Scan | using i_l_orderkey on l2
20118.627

Index Scan | using pk_orders on orders
Filter: o_orderstatus='F'
5666.513

Materialize
24770.922

## Conclusion

TECHNISCHE
UNIVERSITÄT
DARMSTADT

- Derived an asymmetry-aware model from facts and observable behavior
    - Storage properties
    - Algorithm access behavior
- Comparison shows improved performance on application-class benchmark
    - Average speed-up at 48%
    - Individual query speed-up by up to 549%
- Calibration is attracted by strong effects
    - TPC-H features properties not respected in PostgreSQL's optimizer
    - Available degrees of freedom may got abused to compensate deficiencies
- Additional experiments are required to decide upon original hypothesis
    - Using a special load (i.e. a custom micro-benchmark)
    - Focussing on optimization problems in which asymmetry matters explicitly
- Calibration may be a useful tool find optimal configurations in general
    - Typical queries need to be combined in a repeatable benchmark
    - Can cope with any number of variables

# Bibliography

[Knut1973]    Donald E. Knuth.
              *The Art of Computer Programming, Volume III: Sorting and Searching*.
              Addison-Wesley, 1973.

[Kirk1983]    S. Kirkpatrick, C. D. Gelatt, and M. P. Vecchi.
              *Optimization by Simulated Annealing*.
              In *Science*, vol. 220, no. 4598, pages 671–680, 1983.

[YoBC2009]    Terry Yoshii, Christian Black, and Sudip Chahal.
              Solid-state drives in the enterprise: A proof of concept.
              Intel white paper, Intel Corporation, March 2009.