

Where do performance cliffs come from?

Tomas Vondra <tomas.vondra@enterprisedb.com>



Goal(s) of this talk

- discuss one class of performance issues
 - fairly common problem
 - affects cost-based optimization (inherent issue)
- explain why this happens
- maybe give some mitigation hints
 - but no promises, sorry :-)

What is a performance cliff?

- sudden (step) change of performance
- sudden = not proportional to change in "inputs"
- example
 - `SELECT * FROM my_table WHERE column = $1`
 - value "A" matches 1000 rows, query takes 1000 ms
 - value "B" matches 1050 rows, what duration is "expected"?
 - not much more than 1000ms? what if it takes 10000 ms?

Cost vs. Duration

- most databases rely on cost estimates
 - how much "resources" will the plan require (CPU, I/O)
 - assumption: more resources => more time to execute
- cost is ...
 - monotonic and continuous function
 - ... with respect to costing parameters
 - ... selectivity of WHERE condition, number of groups, ...

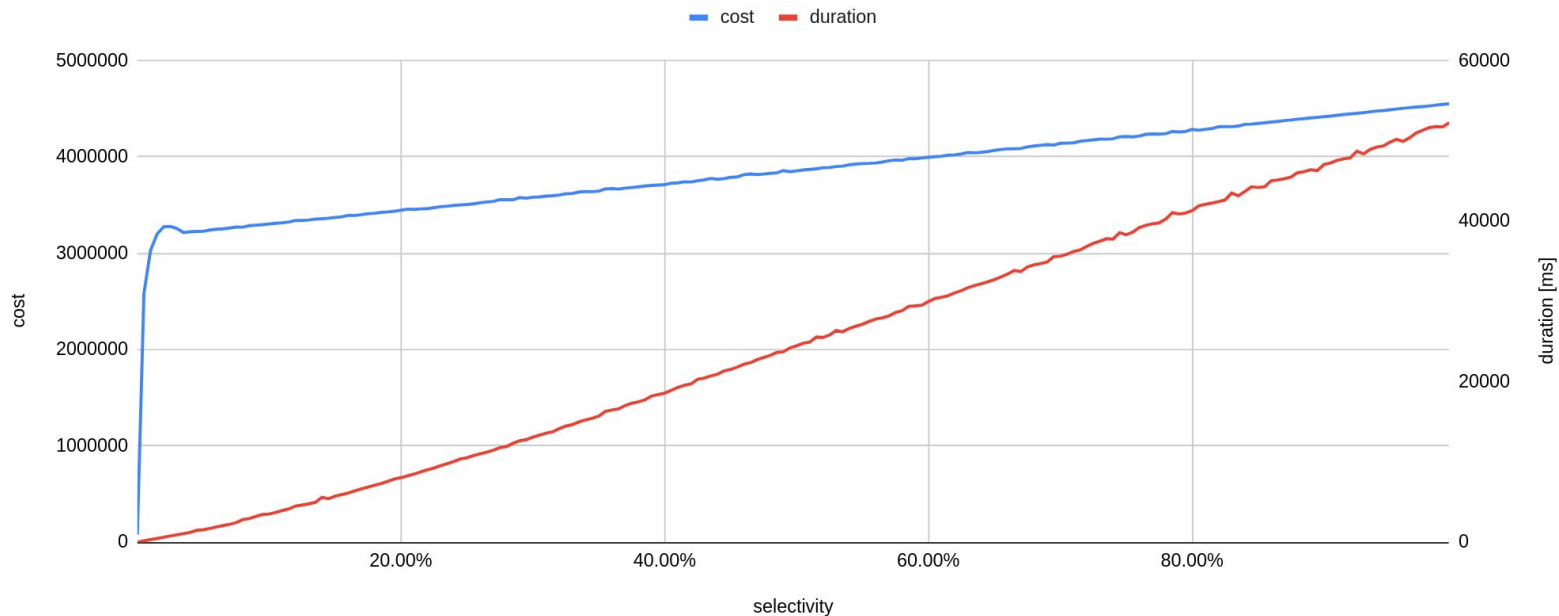
Garbage in - garbage out

- selectivity estimates
- crucial input of the query planning process
- bogus estimate = anything can happen
- we assume selectivities are “good enough”

Example

small selectivity difference => small cost difference => small duration difference

bitmapscan cost vs. duration



Eh?! Where's the discontinuity?

- before: performance cliff is a sudden change in performance
- just now: cost is nice, smooth, without steps, ...

- cost is not timing, but should be correlated
- But why would the timing change in a step?

Ideas?

- ?
- ?
- ?
- ?

Ideas?

- cost is relies on estimates - if wildly wrong, anything can happen
- various things are ultimately decided at runtime
 - e.g. hashjoin / hashagg spilling, on-disk sort, ...
 - on/off decision - one row triggers a lot of work
- we're dealing with multiple plans
 - the whole point of why we calculate costs
 - cost and duration may not "align" perfectly

Runtime decisions

Example: ... IN (list)

```
CREATE TABLE test (a text);
```

```
INSERT INTO test
```

```
SELECT 'aaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaaa' || md5(random()::text)  
FROM generate_series(1,10000000) s(i);
```

```
VACUUM ANALYZE test;
```

```
-- table has ~965MB
```


Example: ... IN (list)

QUERY PLAN

Seq Scan on test (actual rows=0 loops=1)

Filter: (a = ANY ('{aaaaaaaaaaaaaaaaaaaaaaaaaaa..., ...}'::text[]))

Rows Removed by Filter: 10000000

Planning Time: 0.092 ms

Execution Time: 1386.788 ms

(5 rows)

Example: ... IN (list)

- lookup in hash table with ≥ 9 elements
 - fewer elements \Rightarrow linear search
 - but 9 is hard-coded threshold
- ideal threshold depends on cost of comparison
 - specific to data-type and values (e.g. long prefix like here)
 - impossible to know in advance / during execution

Other runtime decisions

- query with in-memory vs. on-disk sort
- query with hashjoin/hashagg in memory vs. spilling to disk
- JIT can be quite expensive & useless
 - enabled depending on total cost of a query
 - ongoing effort to make more granular

Path switch

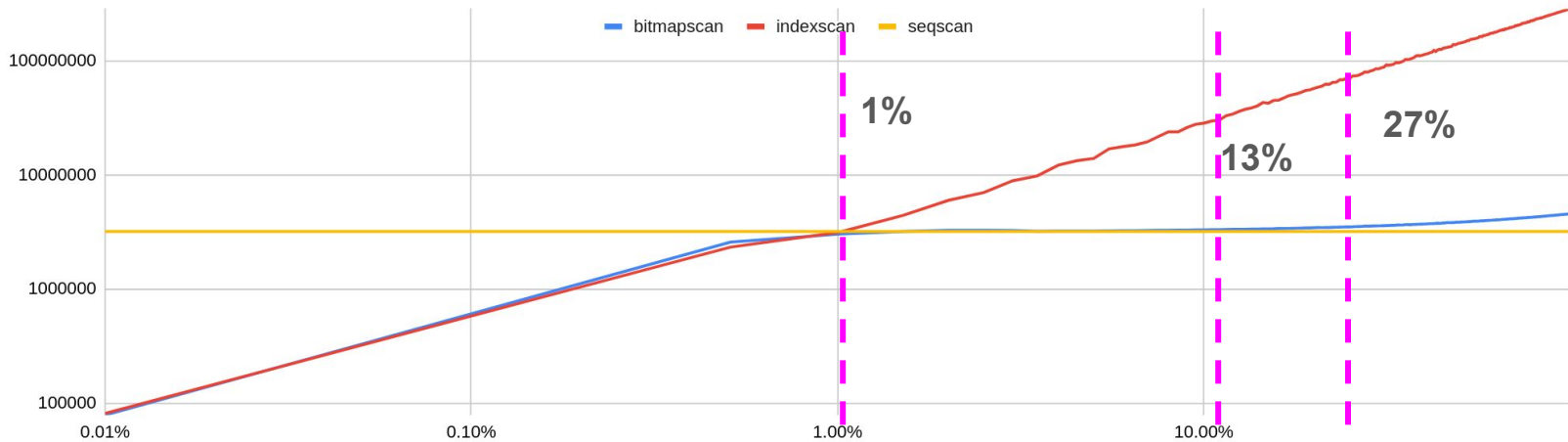
100M rows, random data

```
CREATE TABLE test (a INT, b TEXT) WITH (fillfactor=50);

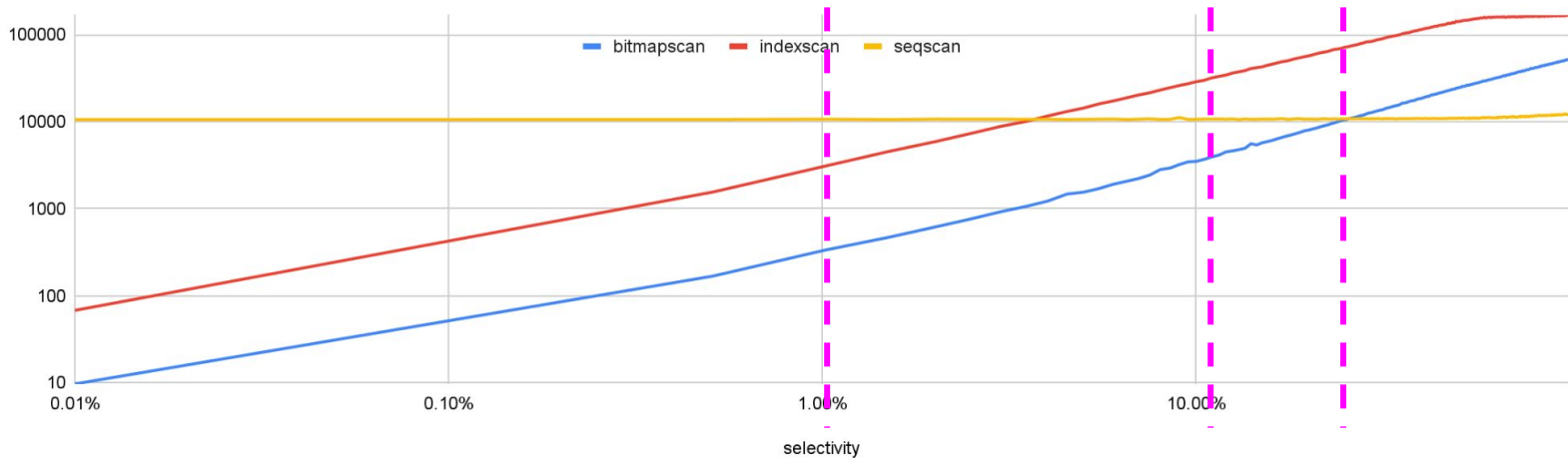
-- 59 rows/page, each page has the same (random) value
INSERT INTO test SELECT a, b FROM (
    SELECT a, b, generate_series(1,59) FROM (
        SELECT 10_000 * random() a,
               md5(random()::text) b
        FROM generate_series(1, 100_000_000/59)
    ) AS x
) AS y;

CREATE INDEX ON test (a);
```

cost: random / 100M rows (SELECT * FROM test WHERE id BETWEEN \$1 AND \$2)



duration: random / 100M rows (SELECT * FROM test WHERE id BETWEEN \$1 AND \$2)





```
SELECT * FROM test WHERE id BETWEEN 1000 AND 1127;  
                QUERY PLAN
```

```
Bitmap Heap Scan on test (actual rows=1293280 loops=1)  
  Recheck Cond: ((id >= 1000) AND (id <= 1127))  
  Heap Blocks: exact=21920  
-> Bitmap Index Scan on test_id_idx (actual rows=1293280 loops=1)  
    Index Cond: ((id >= 1000) AND (id <= 1127))  
Planning Time: 9.268 ms  
Execution Time: 412.993 ms  
(7 rows)
```

```
SELECT * FROM test WHERE id BETWEEN 1000 AND 1128;  
                QUERY PLAN
```

```
Seq Scan on test (actual rows=1301894 loops=1)  
  Filter: ((id >= 1000) AND (id <= 1128))  
  Rows Removed by Filter: 98698091  
Planning Time: 8.289 ms  
Execution Time: 10706.679 ms  
(5 rows)
```

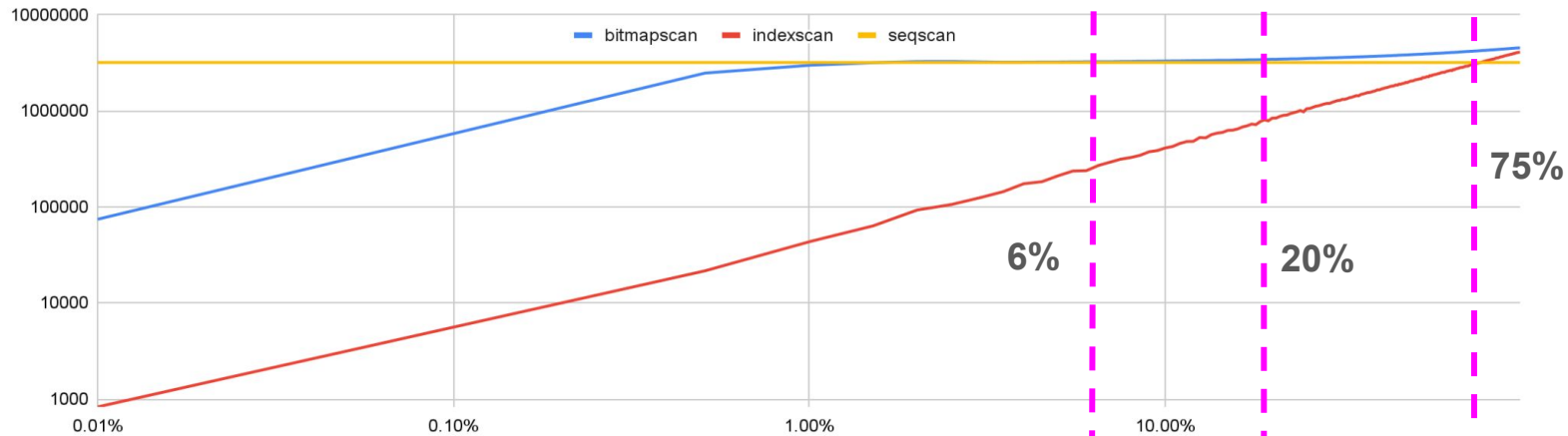
100M rows, sequential/correlated data

```
CREATE TABLE test (a INT, b TEXT) WITH (fillfactor=50);

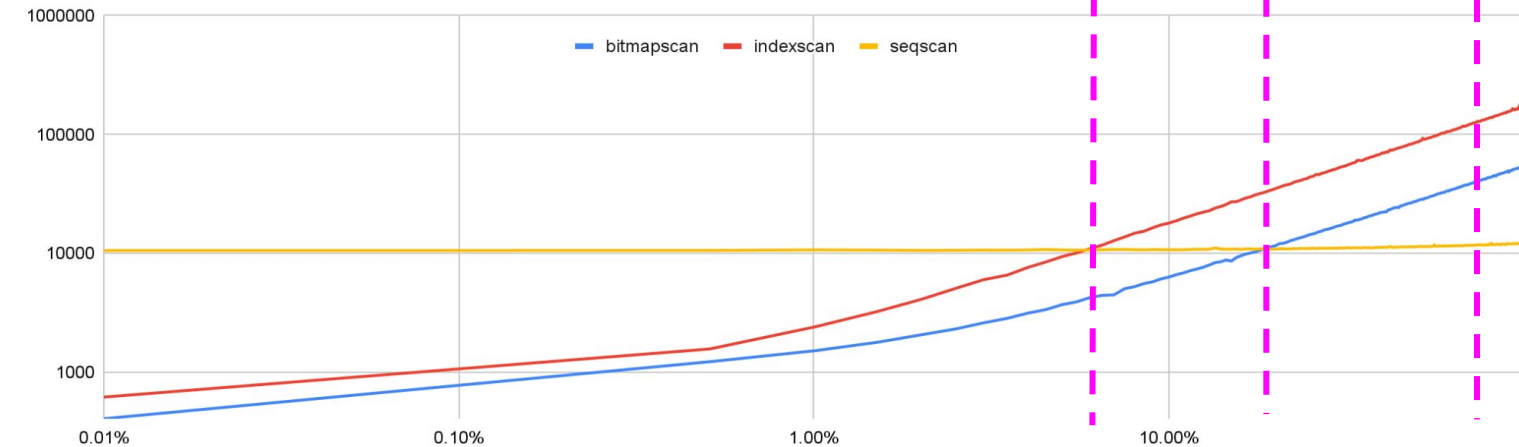
-- monotonic growth, with a bit of random "fuzz"
INSERT INTO test
SELECT (i * 1.0 * 10_000) / 100_000_000 +
       (10_000 * (random() - 0.5)) / 50,
       md5(random()::text)
FROM generate_series(1, 100_000_000) s(i);

CREATE INDEX ON test (a);
```

cost: correlated 100M rows (SELECT * FROM test WHERE id BETWEEN \$1 AND \$2)



duration: correlated 100M rows (SELECT * FROM test WHERE id BETWEEN \$1 AND \$2)





```
select * from test where id between 1000 and 8650;
```

```
QUERY PLAN
```

```
-----  
Seq Scan on test (actual rows=76510346 loops=1)  
  Filter: ((id >= 1000) AND (id <= 8650))  
  Rows Removed by Filter: 23489654  
Planning Time: 0.072 ms  
Execution Time: 11905.432 ms  
(5 rows)
```

```
select * from test where id between 1000 and 8600;
```

```
QUERY PLAN
```

```
-----  
Index Scan using test_id_idx on test (actual rows=76009271 loops=1)  
  Index Cond: ((id >= 1000) AND (id <= 8600))  
Planning Time: 8.398 ms  
Execution Time: 130789.542 ms  
(4 rows)
```

Mitigations?

Mitigations

- really hard to fix (during planning)
- inherent to cost-based planning in general
- costing is approximation
 - simplified model + incomplete data => imperfection
 - G. Graefe: "choice is confusion" [1]
- So, what options do you have?

Mitigations

- try to ensure the "flip" does not trigger
 - increase `work_mem`, for example
 - it "only" moves the threshold ahead
- try to reduce the impact of the "flip"
 - fast but ephemeral storage for temp files?
 - ...

Mitigations

- bit of tuning the cost parameters?
 - `random_page_cost`, `cpu_tuple_cost`, ...
 - can the cost / duration charts align better?
- don't bother to fine-tune the parameter values
 - no parameter value is perfect for all queries
 - the flip needs to happen "close enough"
- some important parameters do not affect costing
 - e.g. `effective_io_concurrency`

Would be better ...

- adaptive execution
 - replace "a priori" decisions with exec time ones
 - ideal: adaptive, smooth transition, not just on/off
 - example: scan type selection vs. "Smooth Scan"
- might also help with estimation errors
- replacement for implementations of a logical node
 - one for scans, another for joins, ...

Robustness / Research papers ...

- Smooth Scan: One Access Path to Rule Them All
R. Borovica, S. Idreos, A. Ailamaki, M. Zukowski, C. Fraser
<https://stratos.seas.harvard.edu/files/stratos/files/smoothscan.pdf>
- A generalized join algorithm
G. Graefe
<https://dl.gi.de/server/api/core/bitstreams/ce8e3fab-0bac-45fc-a6d4-66edaa52d574/content>
- Profile of G. Graefe
https://sigmodrecord.org/publications/sigmodRecord/2009/pdfs/05_Profiles_Graefe.pdf

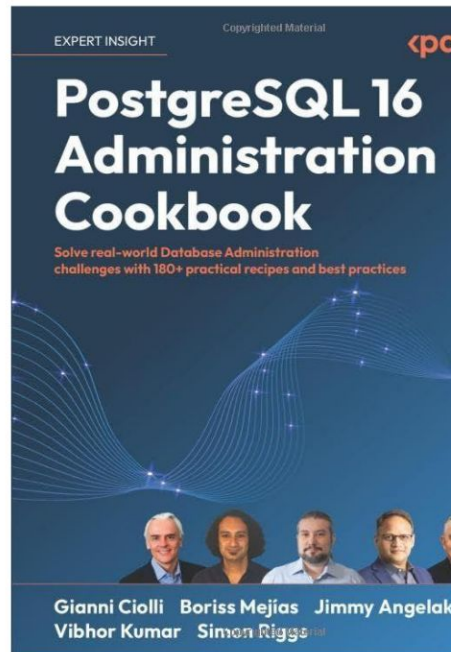


POSTGRESQL 16 ADMINISTRATION COOKBOOK

Take part in EDB's prize draw
to win a brand new book
PostgreSQL 16 Administration Cookbook.

SCAN THE QR CODE AND FILL IN THE FORM TO ENTER

enterprisedb.com



What's a performance cliff?

- sudden (step) change in performance after small change of inputs
- what's an input?
 - not parameter values but rather selectivities of the values
- expectation of "smooth" behavior
 - cost is a continuous function, correlated to duration
 - small change of cost => small change of duration
- what can go wrong?
 - expectation of "sufficiently accurate" estimates => if inputs are bogus, don't expect good plans
 - "smooth cost" applies only to a single path, but we often pick from multiple paths, and the "cost transition points" may not align with the duration (TODO chart comparing cost/duration for scan paths)
 - even a single path may flip between algorithms in slightly inaccurate points (e.g. sort with in-memory vs. on-disk sort or hashagg triggering spill-to-disk), not always known during planning

Examples (single-path)

- IN() clause, with and without hashing (~1000 values?)
- sort with in-memory / on-disk sort
- hash-agg in-memory / spill to disk

Examples (multi-path)

- selecting from multiple scan paths
- cost and duration cross-points may not align
- first show cost chart
- then show duration and how it does not align with cost
- some demos

What can you do?

- not much ;-)
- basic cost tuning to get it "close enough" to duration
- don't skimp on work_mem - if you don't hit the threshold, no cliff
- challenge for optimizer developers
 - every decision = opportunity to get it wrong
 - different algorithm for some parameter values?
 - alternative paths? (new join algorithm, new scan type, ...)
- solution?
 - improve estimates, but don't rely them being 100% correct (literally impossible)
 - focus on "robustness" rather than just raw performance of "ideal plan"
 - adaptive execution - fewer "adaptive" paths rather than many discrete paths
 - examples: SmoothScan and G-join papers (TODO link to papers)
 - examples: maybe unify IndexScan and IndexOnlyScan, make it "gradual"