# Fast faceting using roaring bitmaps

Ants Aasma

www.cybertec-postgresql.com

**CYBERTEC**
DATA SCIENCE & POSTGRESQL

- ► Senior database engineer at Cybertec
- ► I focus on: Performance, high availability, core hacking
- ► We also do: support, HA setups, DBA-as-a-service, private cloud, general consulting

**CYBERTEC**
DATA SCIENCE & POSTGRESQL

► Customer has hundreds of millions of documents.
► Wants to provide faceted search over these documents.
► Need accurate counts, or at worst slightly stale counts.
► Response time below 2 seconds.

# What is faceted search

## Example schema for testing

```
                    Table "test2.documents"
   Column    |           Type          | Collation | Nullable | Default
-------------+-------------------------+-----------+----------+---------
 id          | bigint                  |           | not null |
 created     | timestamp with time zone |          |          |
 finished    | timestamp with time zone |          |          |
 category_id | bigint                  |           | not null |
 tags        | text[]                  |           |          |
 type        | public.mimetype         |           |          |
 size        | bigint                  |           |          |
 title       | text                    |           |          |
Indexes:
    "documents_category_id_idx" btree (category_id)

postgres=# SELECT category_id, ROUDN(COUNT(*)/1e6, 3) millions FROM documents GROUP BY ROLLUP (category_id) ORDER BY 2 DESC LIMIT 6;
 category_id | millions
-------------+----------
             | 100.000
          24 |  60.812
           9 |  15.214
          12 |   6.764
           8 |   3.811
          49 |   2.444
```

**CYBERTEC**
DATA SCIENCE & POSTGRESQL

► Facets:
  ► Created month
  ► Finished month
  ► Category
  ► Type
  ► Size class
► Query: find all other facet counts when filtered by most popular category

▶ Simple option, do a `count(*) group by` for every facet.

```
(SELECT 'created' AS facet_name, date_trunc('month', created)::text AS facet_value,
        COUNT(*) AS cardinality
 FROM documents WHERE category_id = 24 GROUP BY 1, 2)
     UNION ALL
(SELECT 'finished', date_trunc('month', finished)::text, COUNT(*)
 FROM documents  WHERE category_id = 24 GROUP BY 1, 2)
     UNION ALL
(SELECT 'type', type::text, COUNT(*)
 FROM documents WHERE category_id = 24 GROUP BY 1, 2)
     UNION ALL
(SELECT 'size', width_bucket(size, array[0,1000,5000,10000,50000,100000,500000])::text,
        COUNT(*)
 FROM documents WHERE category_id = 24 GROUP BY 1, 2);
```

# Results

38s: https://explain.dalibo.com/plan/99741f1d778c8gcf

**CYBERTEC**
DATA SCIENCE & POSTGRESQL

```
SELECT facet_name, facet_value, COUNT(*) cardinality
FROM documents d, LATERAL (VALUES
    ('created', date_trunc('month', created)::text),
    ('finished', date_trunc('month', finished)::text),
    ('type', type::text),
    ('size', width_bucket(size, array[0,1000,5000,10000,50000,100000,500000
) t(facet_name, facet_value)
WHERE category_id = 24 GROUP BY 1, 2
```

# Results

27s: https://explain.dalibo.com/plan/befh3292ad7cdb12

CYBERTEC
DATA SCIENCE & POSTGRESQL

```
ALTER TABLE documents SET ( parallel_workers = 23 );
```

# More parallelism?

```
ALTER TABLE documents SET ( parallel_workers = 23 );
```

► Time: 18261.594 ms (00:18.262)

**CYBERTEC**
DATA SCIENCE & POSTGRESQL

- ▶ Need to process 60M rows in 2s.
- ▶ For each matching row and each facet:
  - ▶ Calculate the facet value from row
  - ▶ Look up the counter for this facet value and add 1 to it.
- ▶ We have 2s/60e6/4 = 8ns of time to do this
- ▶ RAM access takes ~ 80ns, L3 cache access is 15ns

- ▶ Random access is really bad for CPUs. Can we turn the problem into a sequential one?
- ▶ What if for each facet and value we store a bitmap of matching documents.
- ▶ Counting matches would then be AND'ing together of two bitmaps and counting the number of bits set.
- ▶ CPU's can do this at 128bits*4GHz = 512 Gbit/s/core
- ▶ Could calculate number of matches for 10'000 facet values per second per core.

**CYBERTEC**
DATA SCIENCE & POSTGRESQL

- ▶ Concept is also known as bitmap indexes.
- ▶ Ideally we would want to use different approaches for storing 60M values and 60 values.
- ▶ A popular and fast implementation is Roaring bitmaps. (roaringbitmap.org)
  - ▶ Used by others too including Solr/Elasticsearch, Pinot, Hive, ClickHouse.
- ▶ Good news, there is a postgres extension wrapping it:
  - ▶ github.com/ChenHuajun/pg_roaringbitmap

**CYBERTEC**
DATA SCIENCE & POSTGRESQL

- ► Compressed storage of a set of integers (e.g. document id)
- ► Provides fast SIMD optimized operations for intersection (AND), union (OR) and cardinality (COUNT).
- ► Divides values into 16bit ranges (65'536 values per range) called containers.
- ► First layer is a sorted list of containers that have values, with starting value, type and pointer for each one.
- ► Three types of containers:
  - ► Array of 16bit values (up to 4'096 values)
  - ► 8KB bitmap
  - ► Run length encoded bitmap

**CYBERTEC**
DATA SCIENCE & POSTGRESQL

► Introduces a new roaringbitmap datatype for storing bitmaps in the database
► Some useful functions for us:
   ► rb_build_agg(int) -> roaringbitmap
   ► rb_and(roaringbitmap, roaringbitmap) -> roaringbitmap
   ► rb_or(roaringbitmap, roaringbitmap) -> roaringbitmap
   ► rb_and_agg(roaringbitmap) -> roaringbitmap
   ► rb_cardinality(roaringbitmap) -> int

**CYBERTEC**
DATA SCIENCE & POSTGRESQL

- ► Lets build and maintain a second datastructure of roaring bitmaps for each facet value.
- ► Conceptually an index, but we are doing it in "userspace"
- ► We plan to be modifying data, but mostly only recent rows.
  - ► Store bitmaps as chunks so we can update only modified chunks.
  - ► Capture modifications into a "delta" table and batch apply changes.
- ► Add some functions to make querying easier

# CYBERTEC
DATA SCIENCE & POSTGRESQL

```
CREATE TABLE documents_facets AS
SELECT facet_name, facet_value, id >> 20 chunk_id, rb_build_agg((id & ((1<<20)-1))::int4)
FROM documents d, LATERAL (VALUES
    ('created', date_trunc('month', created)::text),
    ('finished', date_trunc('month', finished)::text),
    ('category_id', category_id::text),
    ('type', type::text),
    ('size', width_bucket(size, array[0,1000,5000,10000,50000,100000,500000])::text)
) t(facet_name, facet_value) GROUP BY 1, 2, 3;

Time: 37334.313 ms (00:37.334)
```

CYBERTEC
DATA SCIENCE & POSTGRESQL

```
CREATE TABLE documents_facets AS
SELECT facet_name, facet_value, id >> 20 chunk_id, rb_build_agg((id & ((1<<20)-1))::int4)
FROM documents d, LATERAL (VALUES
    ('created', date_trunc('month', created)::text),
    ('finished', date_trunc('month', finished)::text),
    ('category_id', category_id::text),
    ('type', type::text),
    ('size', width_bucket(size, array[0,1000,5000,10000,50000,100000,500000])::text)
) t(facet_name, facet_value) GROUP BY 1, 2, 3;

Time: 37334.313 ms (00:37.334)

      table       | table_size | indexes_size
------------------+------------+--------------
 documents        | 20 GB      | 661 MB
 documents_facets | 216 MB     | 848 kB
```
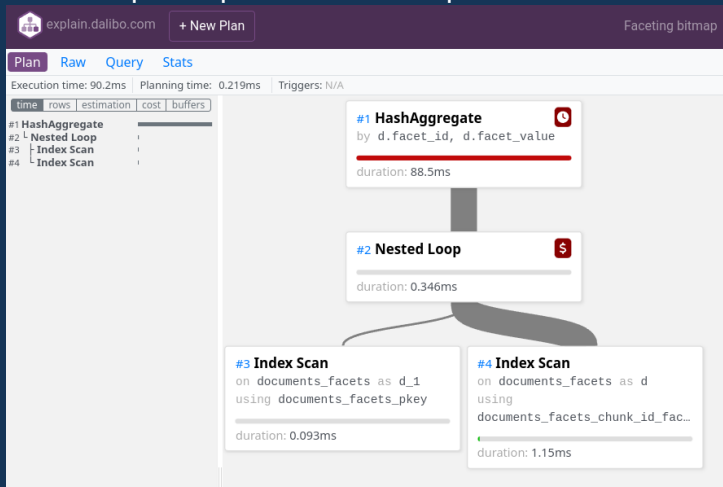
```
WITH lookup AS (
    SELECT chunk_id, postinglist FROM documents_facets d
    WHERE facet_name = 'category_id' AND facet_value = '24'
)
SELECT facet_id, facet_value,
       sum(rb_and_cardinality(lookup.postinglist, d.postinglist)) count
FROM lookup JOIN documents_facets d USING (chunk_id)
WHERE facet_id IN (1,2,4,5)
GROUP BY 1,2;
```

# Results

90ms: https://explain.dalibo.com/plan/71914e60dc7ab324

**CYBERTEC**
DATA SCIENCE & POSTGRESQL

► Don't want to rewrite multiple big bitmaps for each inserted row.
► Create a delta table with structure (facet_id, facet_value, id, delta)
► A trigger on main table that inserts into delta table +1 or -1 for each changed facet.
► Periodically aggregate deltas and merge into main index.

**CYBERTEC**
DATA SCIENCE & POSTGRESQL

- ▶ Started an extension called pgfaceting
  - ▶ github.com/cybertec-postgresql/pgfaceting
- ▶ SQL only, but requires pg_roaringbitmap
- ▶ Open source, help appreciated
- ▶ API is still WIP

# CYBERTEC
DATA SCIENCE & POSTGRESQL

- ▶ Automatic creation and population of facet and delta tables.
- ▶ Facet types:
  - ▶ plain
  - ▶ date_trunc
  - ▶ bucket_facet
- ▶ Delta trigger generation.
- ▶ Query generation facility

```
SELECT faceting.add_faceting_to_table(
    'documents',
    key => 'id',
    facets => array[
        faceting.datetrunc_facet('created', 'month'),
        faceting.datetrunc_facet('finished', 'month'),
        faceting.plain_facet('category_id'),
        faceting.plain_facet('type'),
        faceting.bucket_facet('size',
            buckets => array[0,1000,5000,10000,50000,100000,500000])
    ]
);
-- Add a cron job that does:
CALL faceting.run_maintenance();
```

# Example usage

```
SELECT * FROM faceting.count_results('documents'::regclass,
    filters => array[row('category_id', '24'),
                     row('type', 'image/jpeg')
                ]::faceting.facet_filter[]);
```

**CYBERTEC**
DATA SCIENCE & POSTGRESQL

# Future ideas

- ► Support id's larger than int4
- ► New facet types
  - ► array (e.g. unnest(tags)
  - ► Tree structured facets
  - ► Hierarchical date and time
- ► Automatic range queries
- ► Use delta tables in queries
- ► Full text search

# CYBERTEC
DATA SCIENCE & POSTGRESQL

► Doesn't work well with sparse id spaces like uuid, snowflake.
   Also rules out using ctid
► Tons of unique values (tens to hundreds of thousands)
► Heavily updated facets

# Thanks

- Q & A