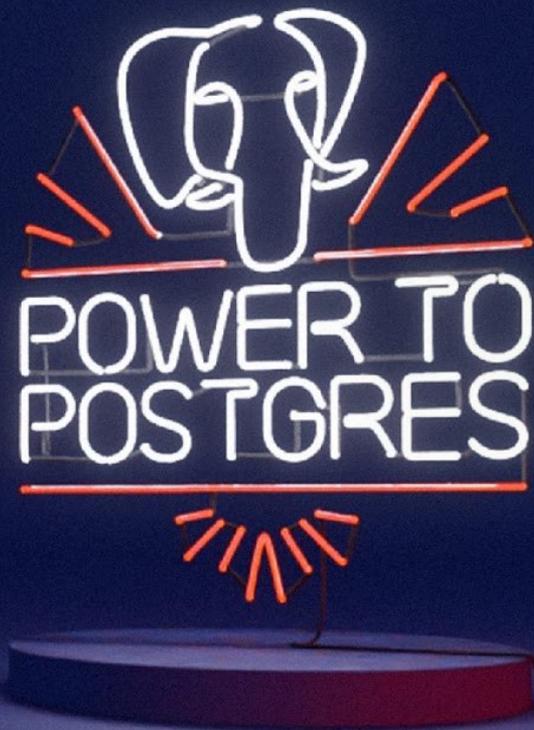


Table Partitioning Transparent but No Magic

Boriss Mejías
Holistic System Software Engineer
Air Guitar Player

22 March 2022 – Nordic PgDay - Helsinki



**It is all about data
a lot of data**

Sadie Jones • Eikä rakkaus oli totta

SADIE JONES • Kutsunottomat vierat

Sadie Jones • Kutsunottomat vierat

Johani Jussi ja Vassu

Maria Jouti Arkielämää

CHUANG-TSE
JOUTI-
LAAN
VAELLUK-
SESTA

ODYSSEUS

Joyce

JAMES JOYCE
ULYSSES

∞

JAMES JOYCE
TAITELIJAN ONNAKOVA
NUORUUDEN VUOSILTA

ODOTTAMATON TOIVOIRETKI

RACHEL JOYCE

HAROLD ERYN

Miranda
July

Katupoika Viipurista

Maria Lehto, Hanna

Maria Lehto, Hanna

TUNTEMATON LAPSI

Tuomas Järvenpää

MIKKA ELÄMÄ VILJATTI

50 SUOMALAISTA TAISTEJA
P. K. KROMIKOSKI
920 Kan

SOTARAAMATTU
USKOMATTOJA
SUOMALAISIA
SOTIHAINEITA
920 Tuo

SUOMALAISEN YHTEISKUNNAN HISTORIA 1400-2000
SUOMALAISTEN SYMBOLIT
Tönnilänneet Tero Halonen & Laura Aho
920 Suo

5 Suomen varhishistorian • AIKAKIRJA
920 Suo

2 Suomen kulttuurihistoria
920 Suo

1 Suomen kulttuurihistoria
TAMMI
920 Suo

SUOMALAISET
JINNOITTEKSI
1720-luvulta
1850-luvulle
920 Suo

RITARIHUONE JA SUOMEN AATEISSUUVUT JOHANNA AMINOFF-WINBERG
920 Rit

JOLO Kvällarna i Helsingfors
920 Ois

MIKKO METSÄMÄKI JA PETTEN HUSUA AKTIVISTIT
920 Met

MESSENIUS SUOMEN RIIMIKRONIKKA
920 Mes

... sarjanat
920 Met





Español
Libros
1º par

Español
Libros
2º par

Ámbito
Internacional
Español

Ámbito
Operativa
Materiales

Rea
Español
Ámbito
Internacional
1º-3

Ámbito
Internacional
2
(seul 1)

Ámbito
Internacional
3
(seul 1)

Ámbito
Internacional
3
(seul 1)

Quintas
Francés
1º-2

Quintas
Francés
1º-2 + 3
+ Anexos

Quintas
Francés
1º-2 + 3
+ Anexos

Esselte
1º-6
1º-6
1

**A large table does not scale
and that's why we want
table partitioning**

SUOMI
Sublime & Ultimate
Online Music Interface

Online music player collecting stats

Register every time a song is played

- User
- Timestamp
- Song
- Artist
- Metadata

Note: Users opt-in to send stats

Online music player collecting stats

```
CREATE TABLE play_counts (  
    user          text  
    , played_on   timestampz  
    , song        text  
    , artist      text  
    , metadata    jsonb  
  
);
```

Online music player collecting stats

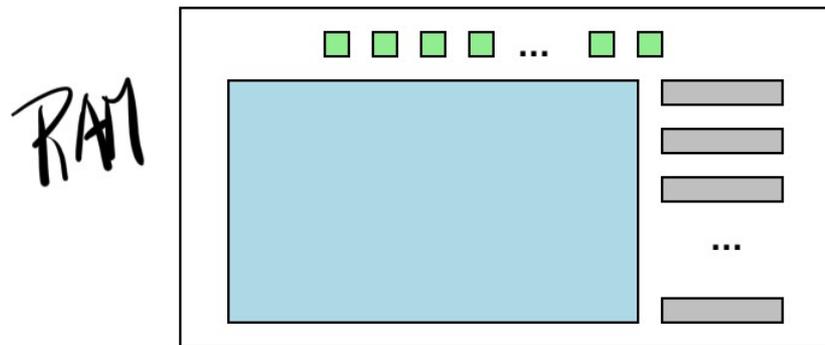
```
CREATE TABLE play_counts (  
    user_id      bigint          REFERENCES users(user_id)  
    , played_on  timestampz  
    , song_id    bigint          REFERENCES songs(song_id)  
    , artist_id  bigint          REFERENCES artists(artist_id)  
    , metadata   jsonb  
  
);
```

Online music player collecting stats

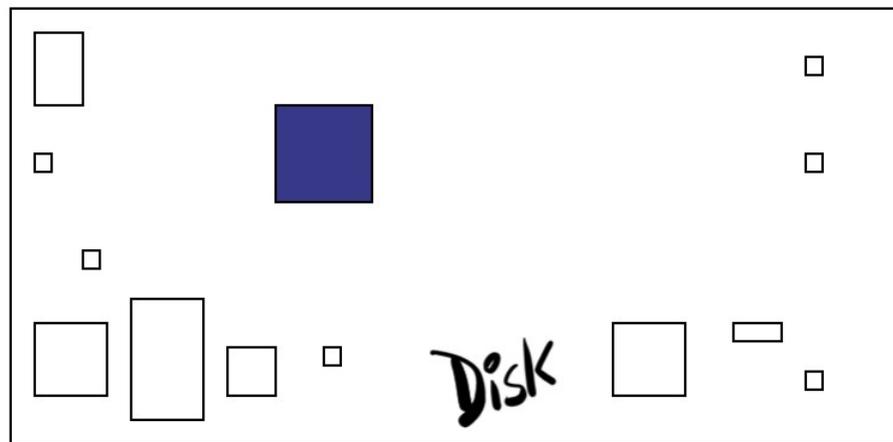
```
CREATE TABLE play_counts (  
    user_id      bigint          REFERENCES users(user_id)  
    , played_on  timestamp(0) WITH TIME ZONE  
    , song_id    bigint          REFERENCES songs(song_id)  
    , artist_id  bigint          REFERENCES artists(artist_id)  
    , metadata   jsonb  
  
);
```

Online music player collecting stats

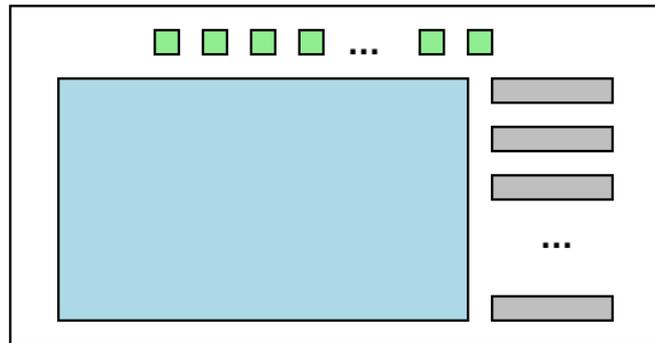
```
CREATE TABLE play_counts (  
    user_id      bigint          REFERENCES users(user_id)  
    , played_on  timestamp(0) WITH TIME ZONE  
    , song_id    bigint          REFERENCES songs(song_id)  
    , artist_id  bigint          REFERENCES artists(artist_id)  
    , metadata   jsonb  
    , PRIMARY KEY (user_id, played_on)  
);
```



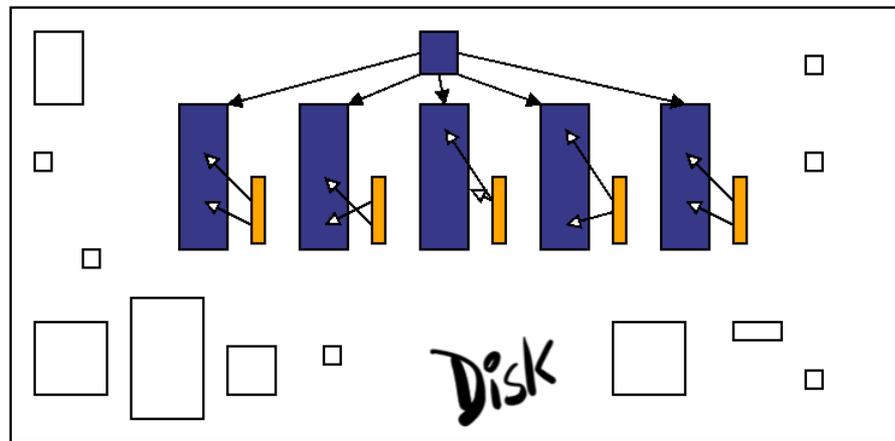
↕ I/O



RAM



↕ I/O

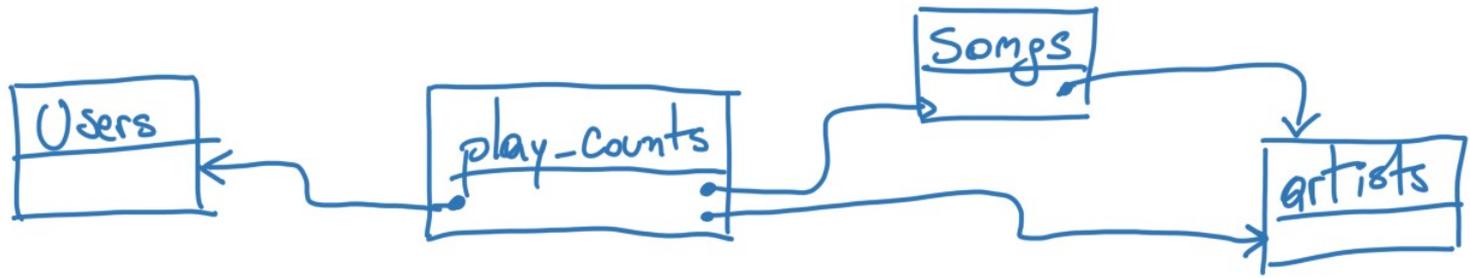


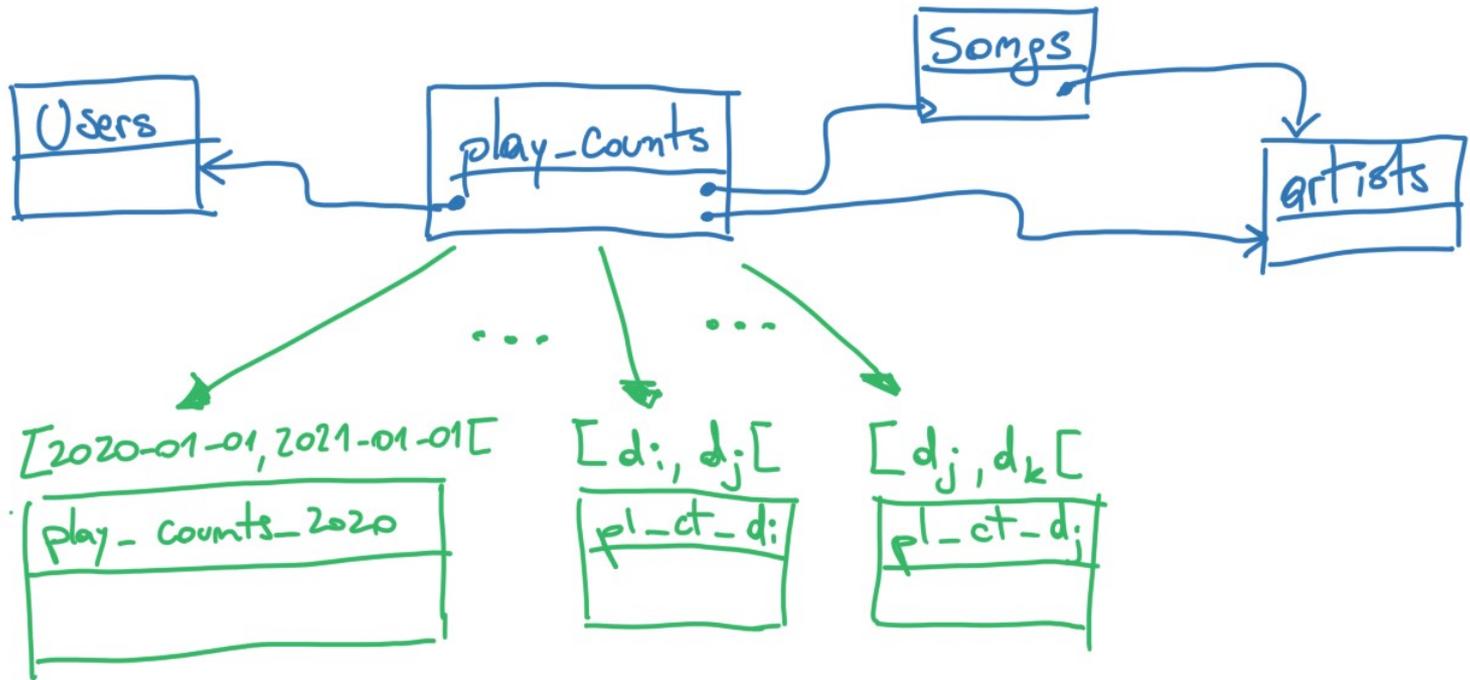
Online music player collecting stats

```
CREATE TABLE play_counts (  
    user_id      bigint          REFERENCES users(user_id)  
    , played_on  timestamp(0) WITH TIME ZONE  
    , song_id    bigint          REFERENCES songs(song_id)  
    , artist_id  bigint          REFERENCES artists(artist_id)  
    , metadata   jsonb  
    , PRIMARY KEY (user_id, played_on)  
);
```

Online music player collecting stats

```
CREATE TABLE play_counts (  
    user_id      bigint          REFERENCES users(user_id)  
    , played_on  timestamp(0) WITH TIME ZONE  
    , song_id    bigint          REFERENCES songs(song_id)  
    , artist_id  bigint          REFERENCES artists(artist_id)  
    , metadata   jsonb  
    , PRIMARY KEY (user_id, played_on)  
)  
PARTITION BY RANGE (played_on);
```





Online music player collecting stats

```
CREATE TABLE play_counts (  
    user_id      bigint          REFERENCES users(user_id)  
    , played_on  timestamp(0) WITH TIME ZONE  
    , song_id    bigint          REFERENCES songs(song_id)  
    , artist_id bigint          REFERENCES artists(artist_id)  
    , metadata   jsonb  
    , PRIMARY KEY (user_id, played_on)  
) PARTITION BY RANGE (played_on);
```

Creating partitions

```
CREATE TABLE play_counts_2020 PARTITION OF play_counts  
FOR VALUES FROM ('2020-01-01') TO ('2021-01-01');
```

Creating partitions

```
CREATE TABLE play_counts_2020 PARTITION OF play_counts  
FOR VALUES FROM ('2020-01-01') TO ('2021-01-01');
```

```
CREATE TABLE play_counts_2021q1 PARTITION OF play_counts  
FOR VALUES FROM ('2021-01-01') TO ('2021-04-01');
```

Creating partitions

```
CREATE TABLE play_counts_2020 PARTITION OF play_counts  
FOR VALUES FROM ('2020-01-01') TO ('2021-01-01');
```

```
CREATE TABLE play_counts_2021q1 PARTITION OF play_counts  
FOR VALUES FROM ('2021-01-01') TO ('2021-04-01');
```

```
CREATE TABLE play_counts_2021q2 PARTITION OF play_counts  
FOR VALUES FROM ('2021-04-01') TO ('2021-07-01');
```

```
CREATE TABLE play_counts_2021q3 PARTITION OF play_counts  
FOR VALUES FROM ('2021-07-01') TO ('2021-10-01');
```

```
CREATE TABLE play_counts_2021q4 PARTITION OF play_counts  
FOR VALUES FROM ('2021-10-01') TO ('2022-01-01');
```

Let's see it working

Understanding what is going on

```
EXPLAIN (ANALYZE, BUFFERS)
SELECT COUNT(*)
FROM play_counts;
```

Understanding what is going on

```
EXPLAIN (ANALYZE, BUFFERS)
SELECT COUNT(*)
FROM play_counts
WHERE played_on >= now() - '1 month'::interval;
```

**Postgres can prune tables
based on the partition key**

played_on

**What about queries
based on the user?**

Query across all partitions

```
EXPLAIN ANALYZE  
SELECT *  
FROM play_counts  
WHERE user_id = 10;
```

Dynamically prune partitions

```
EXPLAIN ANALYZE
SELECT *
FROM play_counts
WHERE user_id = 10
ORDER BY played_on DESC
LIMIT 200;
```

Getting the Query right is Fundamental

From

```
SELECT COUNT(*)  
FROM play_counts p  
JOIN artists a ON p.artist = a.artist_id  
WHERE a.name='Korpiklaani';
```

To

```
SELECT COUNT(*)  
FROM play_counts p  
WHERE artist = 22780  
AND played_on >= '2021-11-21';
```

Consider CTEs – WITH Queries

```
WITH artist AS (  
    SELECT artist_id, added_on  
    FROM artists WHERE name = 'Korpiklaani'  
)  
SELECT COUNT(*)  
FROM play_counts  
WHERE artist = (SELECT artist_id FROM artist)  
AND played_on >= (SELECT added_on FROM artist);
```

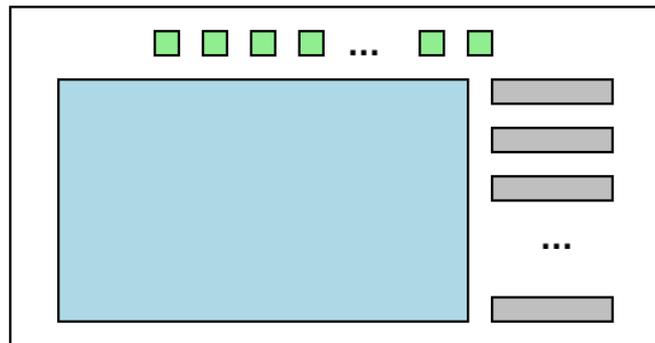
What about indexes?

Create index at the parent table

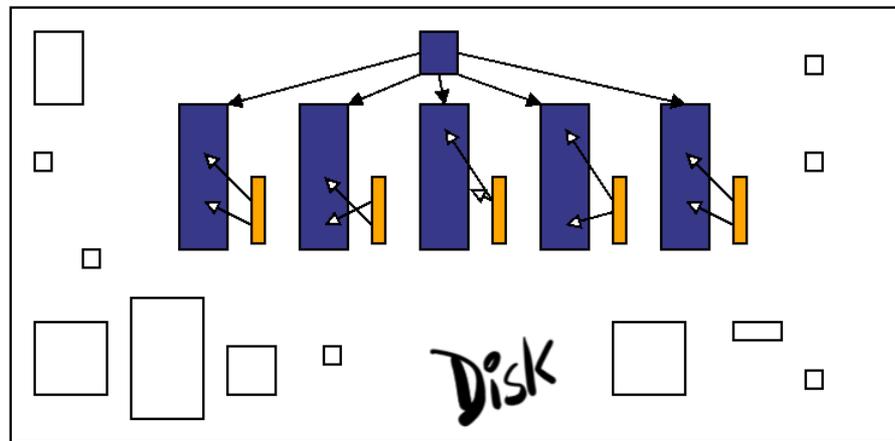
```
CREATE INDEX play_counts_artist_idx  
ON play_counts (artist_id);
```

NOTE: This is NOT a global index!

RAM



↕ I/O



**Before we
continue**



**Partition large tables
when they have a logical
partition key**

**Make sure your queries match
the way tables are partitioned**

More partition strategies



50 SUOMALAISTA TAISTEJA
P. K. KROMIKOSKI
920 Kan

SOTARAAMATTU
USKOMATTOJA
SUOMALAISIA
SOTIHAINEITA
920 Tuo

SUOMALAISEN YHTEISKUNNAN HISTORIA 1400-2000
SUOMALAISTEN SYMBOLIT
Tönnilänter, Tero Halonen & Laura Aho
920 Suo

5 Suomen varhishistorian • AIKAKIRJA
920 Suo

2 Suomen kulttuurihistoria
920 Suo

1 Suomen kulttuurihistoria
920 Suo

SUOMALAISET
JINNOITTEKSI
1720-luvulta
1850-luvulle
920 Suo

RITARIHUONE JA SUOMEN AATEISSUUVUT JOHANNA AMINOFF-WINBERG
920 Rit

JOLO Kvällarna i Helsingfors
920 Ois

MIKKO METSÄMÄKI JA PETTERI HUSUA AKTIVISTIT
920 Met

MESSENIUS SUOMEN RIIMIKRONIKKA
920 Mes

... sarjanat
920 Met



Table partitioning by HASH value

```
CREATE TABLE artists (  
    artist_id    bigint GENERATED ALWAYS AS IDENTITY  
    , name       text  
    , added_on   date  
    , axed_on    date  
    , PRIMARY KEY (artist_id)  
) PARTITION BY HASH (artist_id);
```

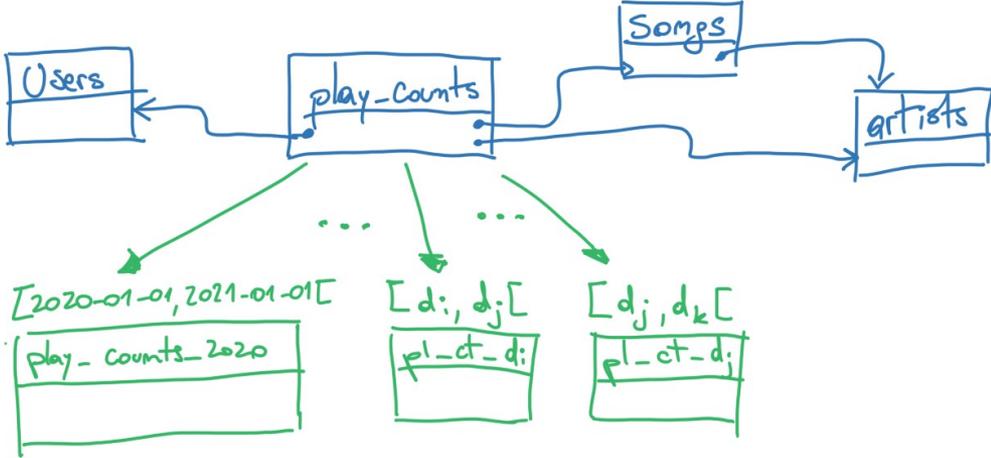
Table partitioning by HASH value

```
CREATE TABLE artists_0 PARTITION OF artists  
FOR VALUES WITH (modulus 3, remainder 0);
```

```
CREATE TABLE artists_1 PARTITION OF artists  
FOR VALUES WITH (modulus 3, remainder 1);
```

```
CREATE TABLE artists_2 PARTITION OF artists  
FOR VALUES WITH (modulus 3, remainder 2);
```

Complete support for Foreign Keys



Complete support for Foreign Keys

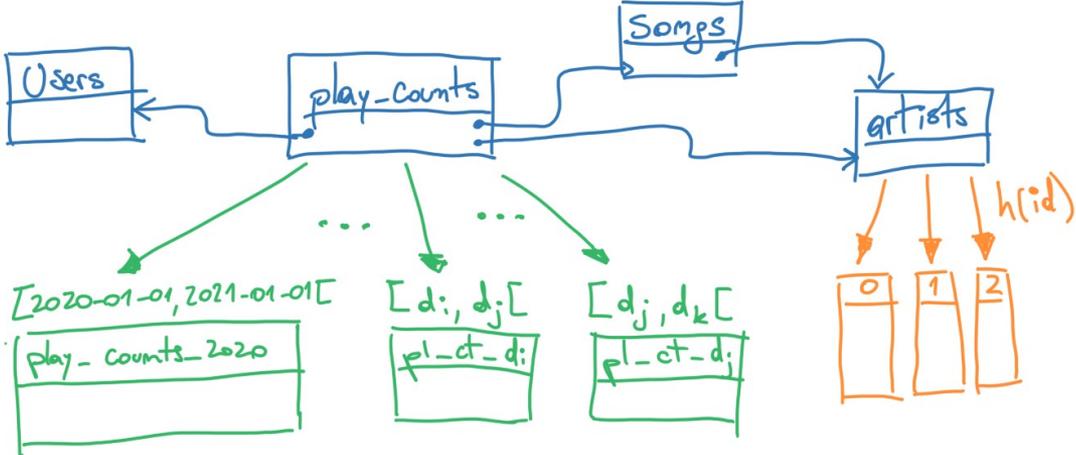


Table partitioning by LIST

```
CREATE TABLE playlists (  
    id          uuid  
    , name      text  
    , category  text  
    , tracks    jsonb  
) PARTITION BY LIST (category);
```

Table partitioning by LIST

```
CREATE TABLE playlists_kids PARTITION OF playlists  
FOR VALUES IN ('kids');
```

```
CREATE TABLE playlists_16 PARTITION OF playlists  
FOR VALUES IN ('16+', 'explicit', 'heartbreaker');
```

The DEFAULT Partition

```
CREATE TABLE playlists_kids PARTITION OF playlists  
FOR VALUES IN ('kids');
```

```
CREATE TABLE playlists_16 PARTITION OF playlists  
FOR VALUES IN ('16+', 'explicit', 'heartbreaker');
```

```
CREATE TABLE playlists_default PARTITION OF playlists  
DEFAULT;
```

Inserting data

```
INSERT INTO playlists
VALUES (gen_random_uuid(), 'boomhut', 'kids', NULL);
-- goes into playlists_kids
```

Inserting and Updating data

```
INSERT INTO playlists
VALUES (gen_random_uuid(), 'boomhut', 'kids', NULL);
-- goes into playlists_kids
```

```
UPDATE playlists
SET category = 'explicit' WHERE name='boomhut';
-- from playlists_kids to playlists_16
```

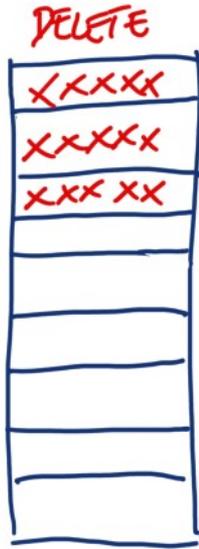
Inserting more data

```
INSERT INTO playlists  
VALUES (gen_random_uuid(), 'la plage', 'senior', NULL);  
-- goes into playlists_default
```

Maintenance

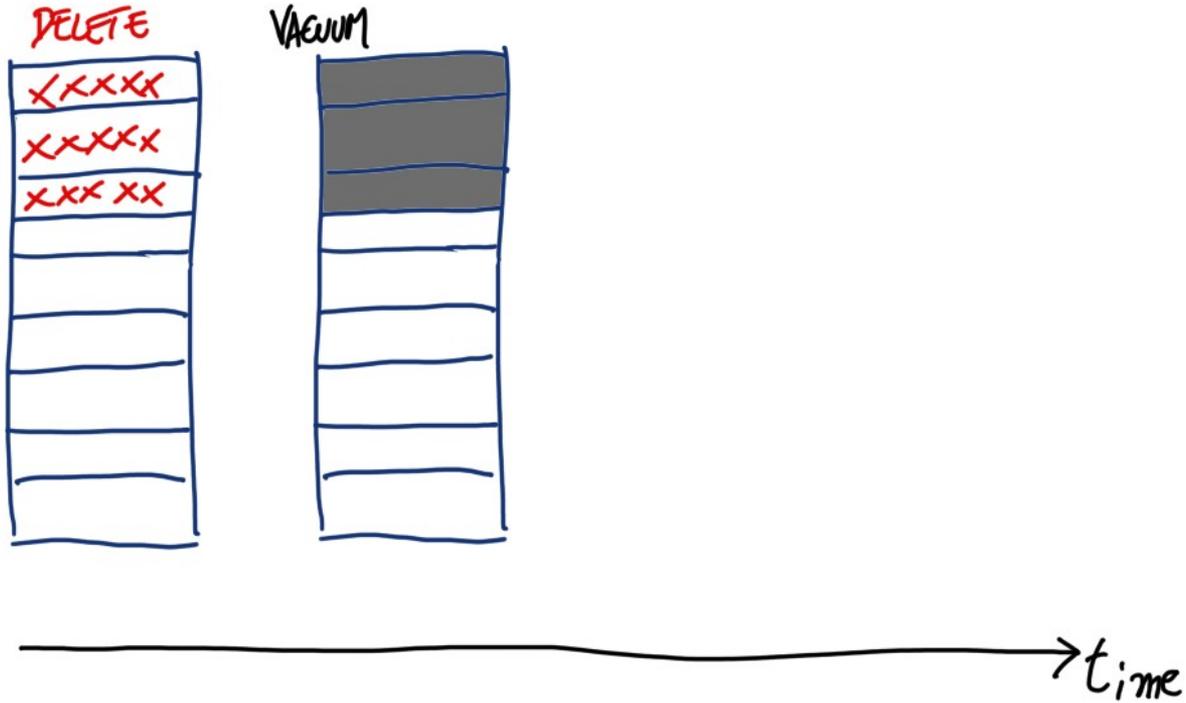


Delete data in Postgres

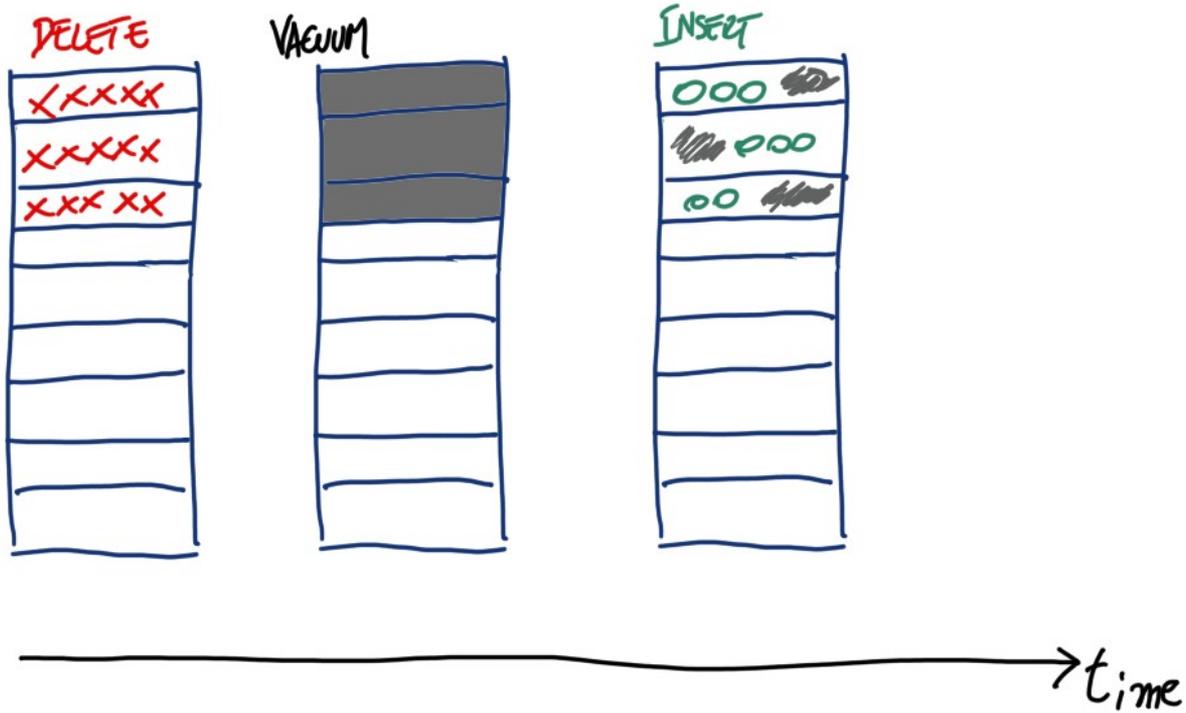


→ time

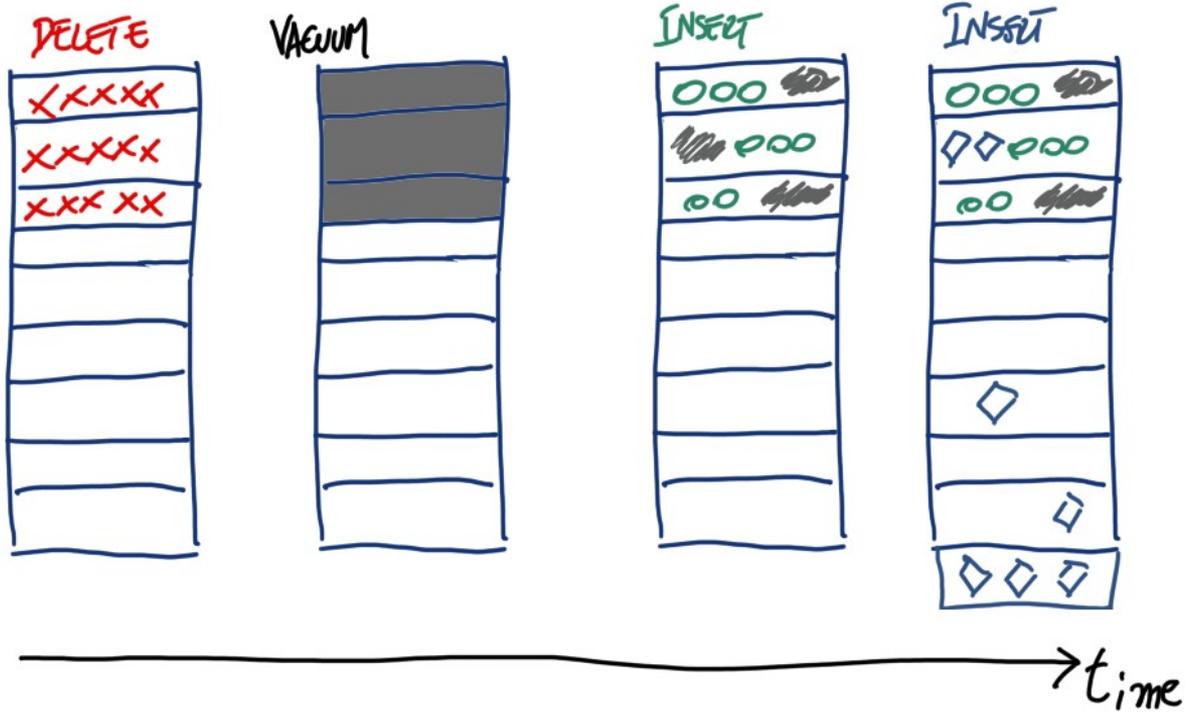
Delete data in Postgres



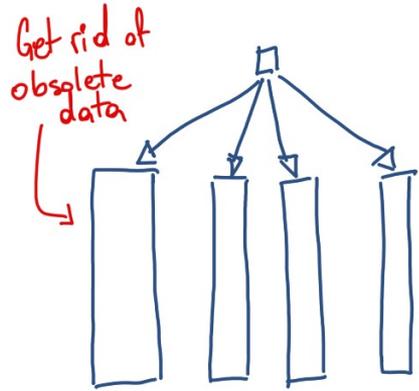
Delete data in Postgres



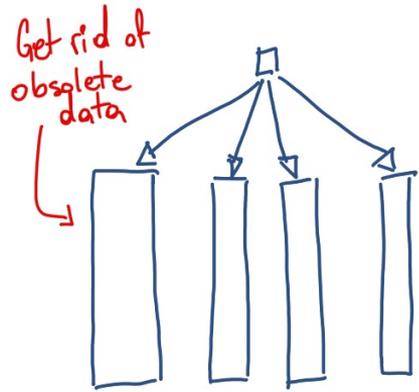
Delete data in Postgres



Detach or Drop partition

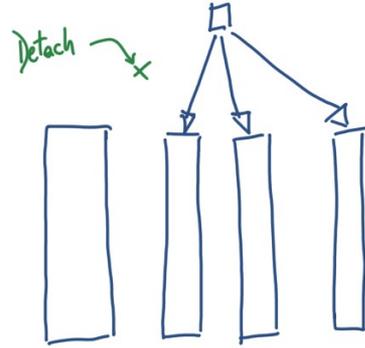


Detach or Drop partition

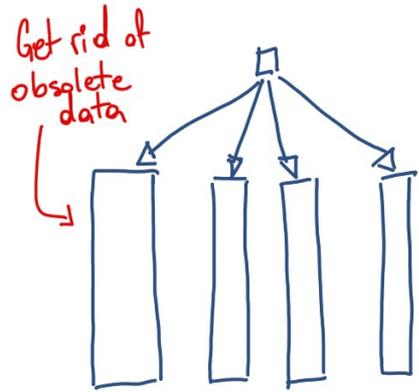


Detach to archive

A green arrow points from the text "Detach to archive" towards the right diagram.

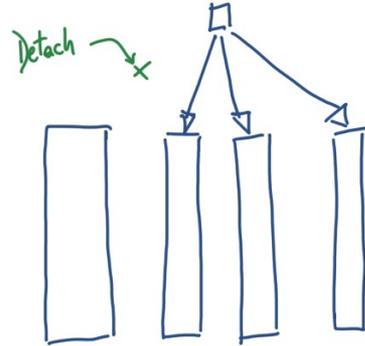


Detach or Drop partition

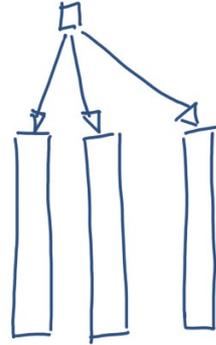


Detach to archive

just drop it



~~DROP TABLE~~



Detach a partition

```
ALTER TABLE play_counts DETACH PARTITION play_counts_2019;
```

Dropping a partition is dropping a table

```
DROP TABLE play_counts_2019;
```

Getting rid of obsolete data

```
ALTER TABLE play_counts DETACH PARTITION play_counts_2019;
```

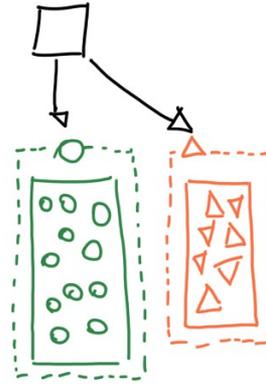
```
DROP TABLE play_counts_2019;
```

Create new partitions

```
CREATE TABLE play_counts_202204 PARTITION OF play_counts  
FOR VALUES FROM ('2022-04-01') TO ('2022-05-01');
```

**What about importing
a full partition?
What about bulk uploads?**

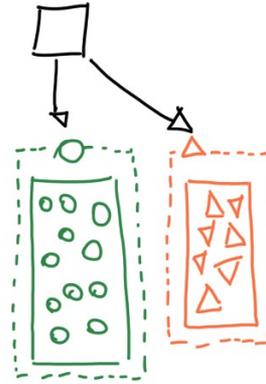
Attach partition with data



→ time

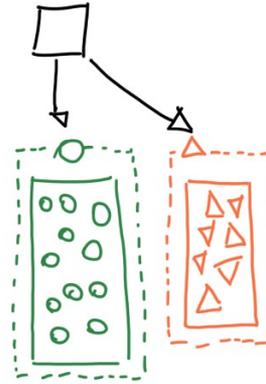
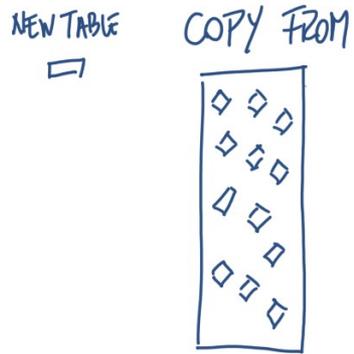
Attach partition with data

NEW TABLE
□



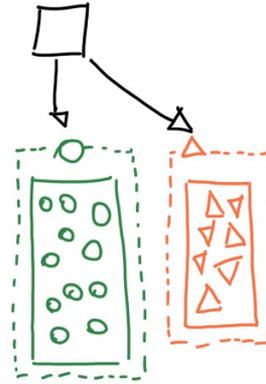
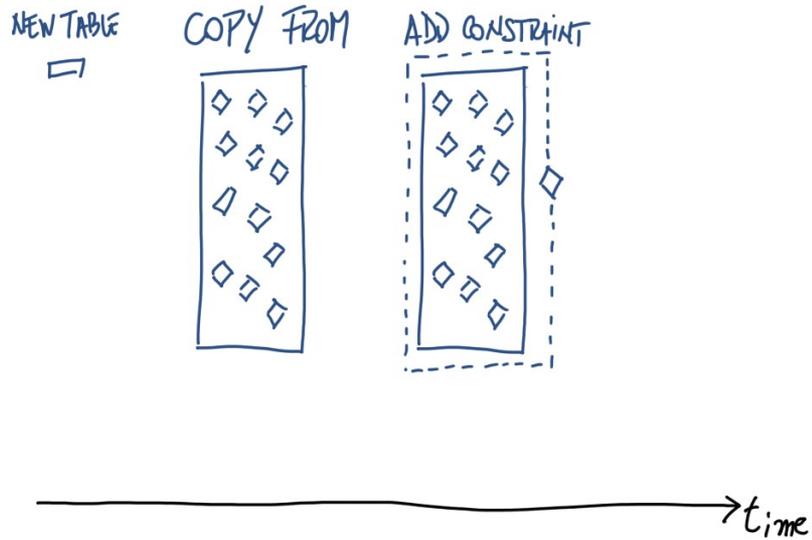
→ time

Attach partition with data

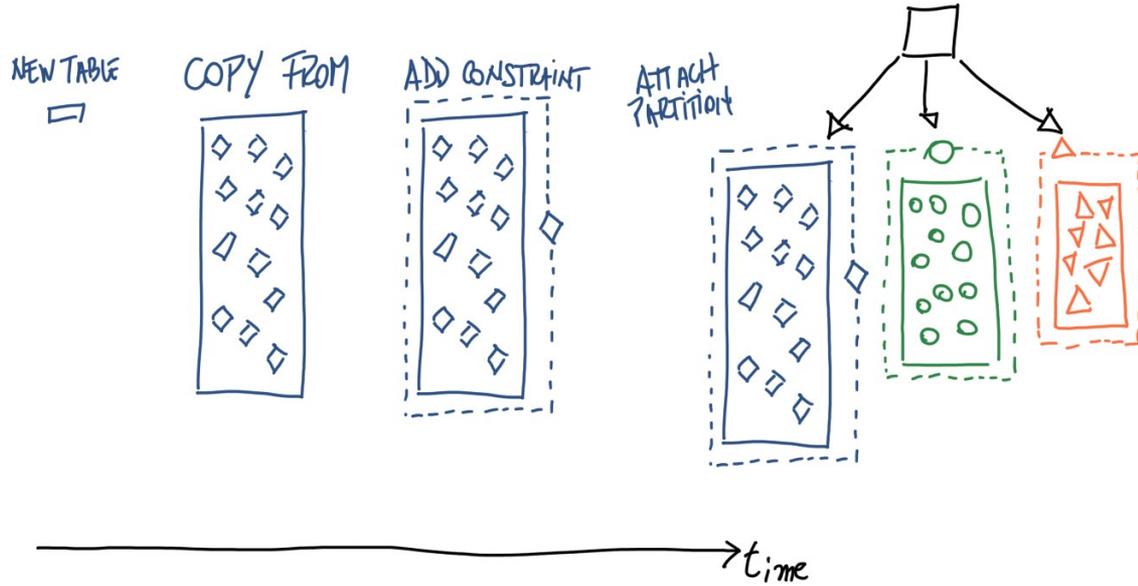


→ time

Attach partition with data

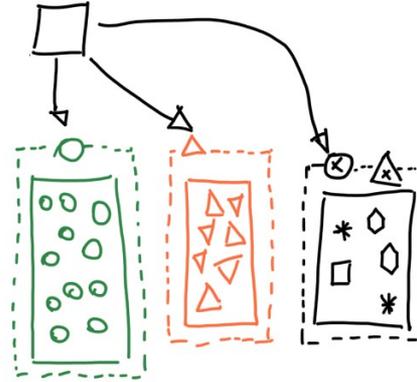


Attach partition with data



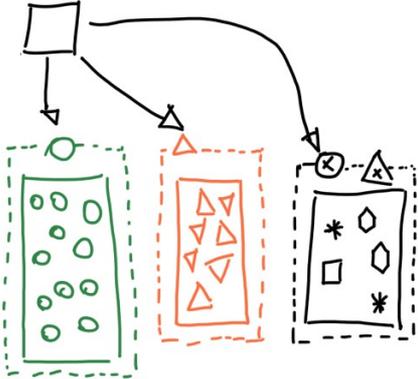
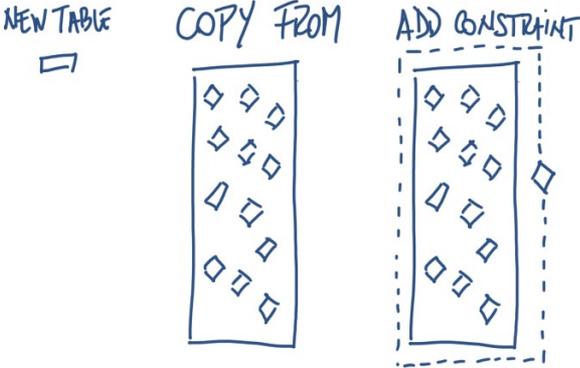
**Remember
No Magic**

The price of Default partitions



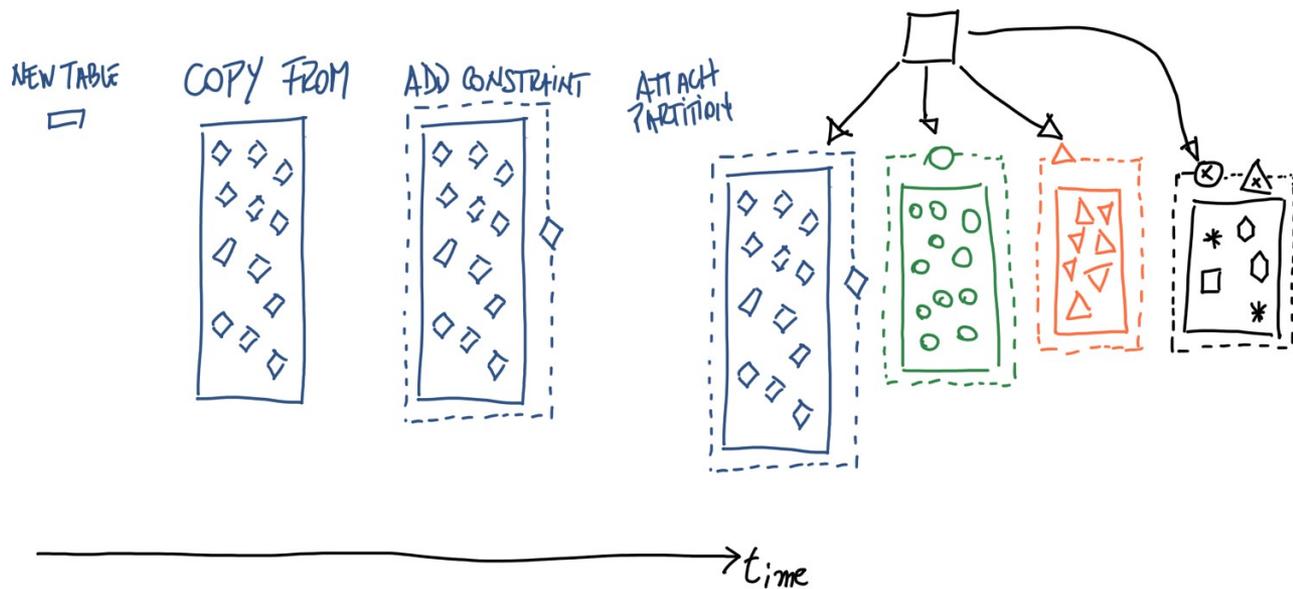
→ time

The price of Default partitions

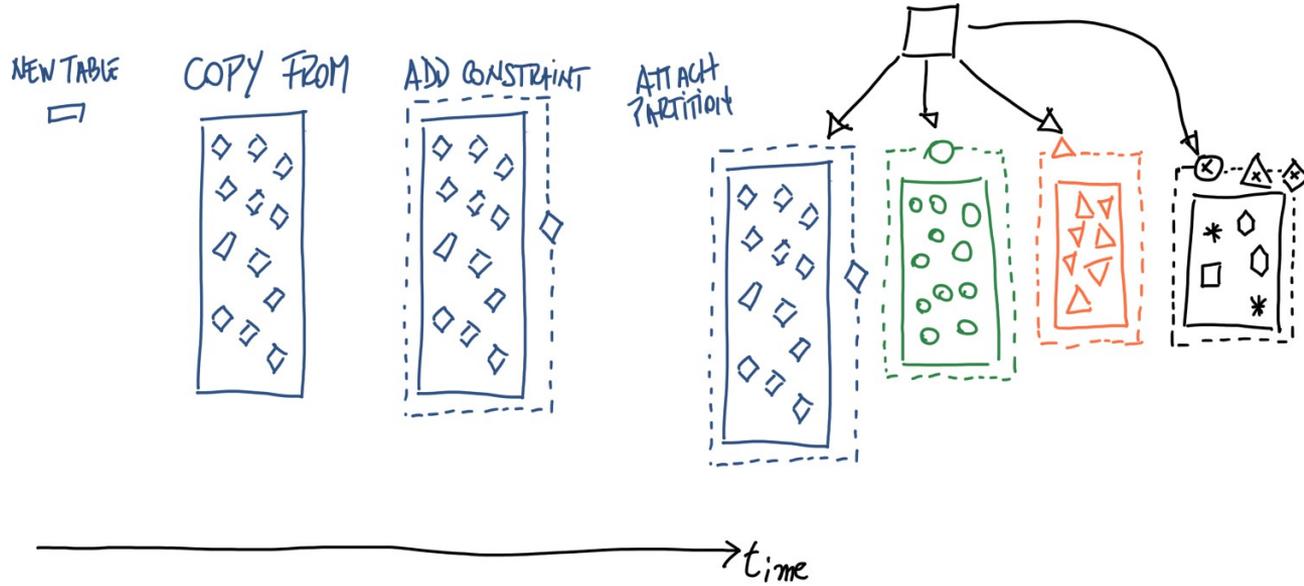


→ time

The price of Default partitions



The price of Default partitions



Import new partition

```
CREATE TABLE playlists_headbanging  
  (LIKE playlists including all);
```

```
\COPY playlists_headbanging  
FROM '/tmp/headbanging_postgres.txt' CSV DELIMITER ',';
```

Add Constraint before attaching

```
ALTER TABLE playlist_headbang
ADD CONSTRAINT headbanger
CHECK ((category IS NOT NULL) AND
       (category = ANY (ARRAY['rock', 'metal'])));
```

And attach

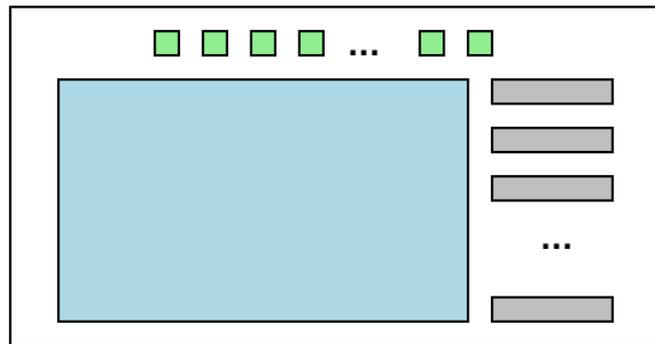
```
ALTER TABLE playlist_headbang
ADD CONSTRAINT headbanger
CHECK ((category IS NOT NULL) AND
       (category = ANY (ARRAY['rock', 'metal'])));
```

```
ALTER TABLE playlists
ATTACH PARTITION playlist_headbang
FOR VALUES IN ('metal', 'rock');
```

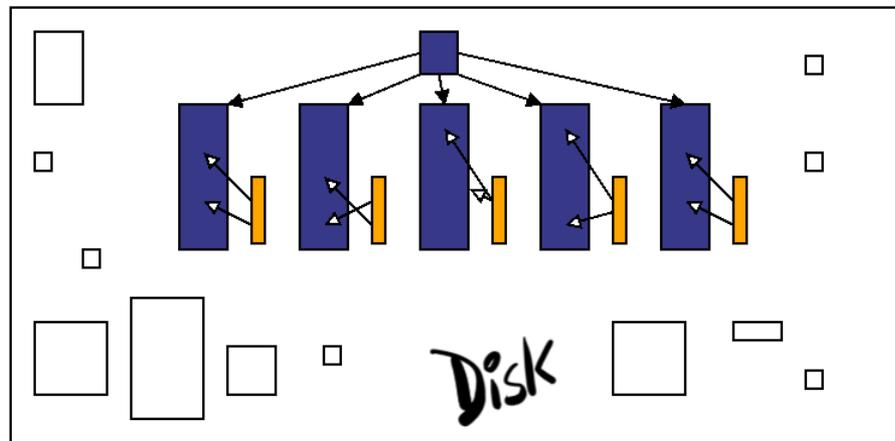
Declarative Table Partitioning

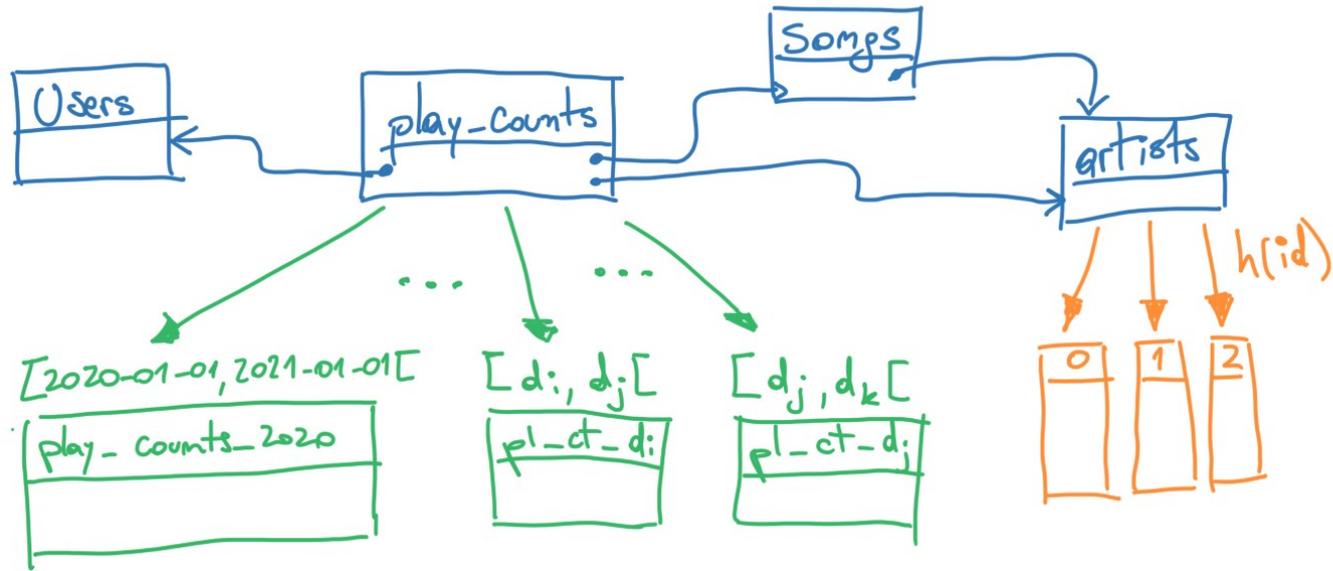


RAM

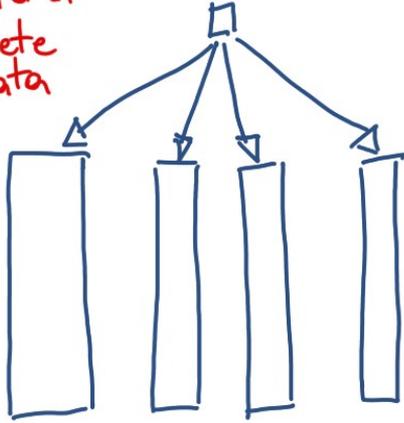


↕ I/O





Get rid of
obsolete
data



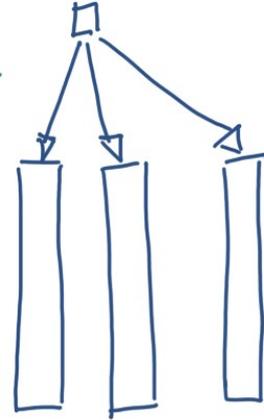
Detach to
archive



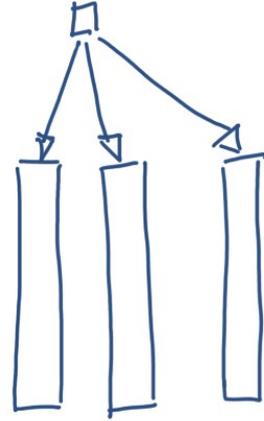
just
drop it



Detach



~~DROP TABLE~~

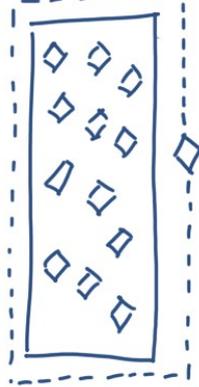


NEW TABLE
□

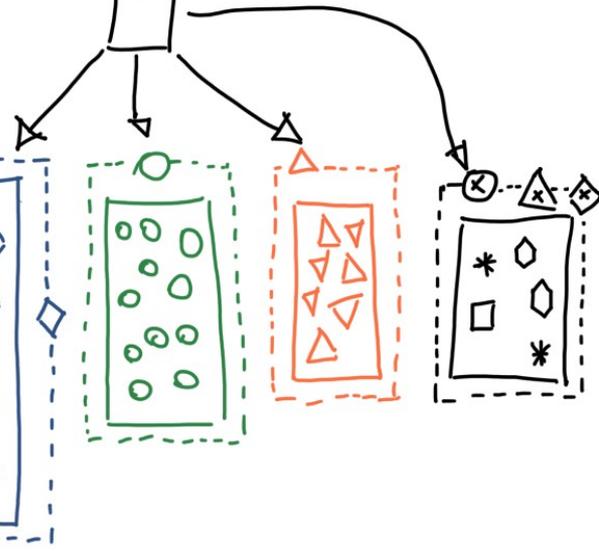
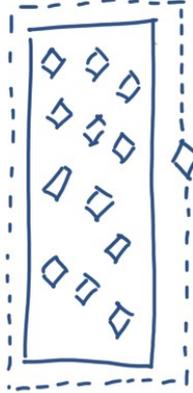
COPY FROM



ADD CONSTRAINT



ATTACH PARTITION



→ time

**And don't forget to
EXPLAIN and ANALYZE
your queries.
You may need to redesign.**

Table Partitioning Transparent but No Magic

Kiitos

boriss.mejias@enterprisedb.com
@tchorix

